

Properties of neutrality tests based on allele frequency spectrum

L. Ferretti^{a,b,*}, G. Marmorini^c, S. Ramos-Onsins^{a,b}

^a*Department of Animal Science and Food, Facultat de Veterinària, Universitat Autònoma de Barcelona, 08193 Bellaterra, Spain*

^b*Centre for Research in Agricultural Genomics (CRAG), 08193 Bellaterra, Spain*

^c*Department of Physics, Keio University, 223-8521 Kanagawa, Yokohama, Hiyoshi 4-1-1, Japan*

Abstract

One of the main necessities for population geneticists is the availability of statistical tools that enable to accept or reject the neutral Wright-Fisher model with high power. A number of statistical tests have been developed to detect specific deviations from the null frequency spectrum in different directions (i.e., Tajima's D, Fu and Li's F and D test, Fay and Wu's H). Recently, a general framework was proposed to generate all neutrality tests that are linear functions of the frequency spectrum. In this framework, a family of optimal tests was developed to have almost maximum power against a specific alternative evolutionary scenario. Following these developments, in this paper we provide a thorough discussion of linear and nonlinear neutrality tests. First, we present the general framework for linear tests and emphasize the importance of the property of scalability with the sample size (that is, the results of the tests should not depend on the sample size), which, if missing, can guide to errors in data interpretation. The motivation and structure of linear optimal tests are discussed. In a further generalization, we develop a general framework for nonlinear neutrality tests and we derive nonlinear optimal tests for polynomials of any degree in the frequency spectrum.

Keywords: Coalescent theory, Site frequency spectrum, Population genetics, Statistical power, Summary statistics

*Corresponding author

Email address: `luca.ferretti@uab.cat` (L. Ferretti)

Contents

1	Introduction	2
2	Linear neutrality tests	4
2.1	General framework	4
2.2	Transformations of weights and invariance of tests	6
2.3	Generalized D' tests for multilocus analysis	7
3	Sample size independent tests	7
3.1	Scaling of weights with sample size	7
3.2	Divergent weights	11
3.3	Weights of singletons	13
3.4	Scaling of weights in tests without an outgroup	15
3.5	Alternative choices of scaling	16
4	Linear optimal tests	17
4.1	On the existence of generic tests	17
4.2	Optimal tests and their geometric structure	18
4.3	Scaling of optimal tests	20
4.4	Generalizing D' for optimal test	21
4.5	Linear tests with maximum power	22
5	Beyond linear neutrality tests	23
5.1	Quadratic and nonlinear tests	23
5.2	Strongly centered optimal tests	27
5.3	Weakly centered optimal tests	29
5.4	Simulations of the power of optimal tests	31
6	Conclusions	33
A	Moments of the frequency spectrum in the independent sites approximation	35
B	Proofs	38

1. Introduction

Statistical tests for neutrality are important and useful tools for population genetics. Since the development of molecular genetics techniques allowed

to obtain nucleotide sequences for the study of populations genetics [1], a number of neutrality tests have been developed with the objective to facilitate the interpretation of an increasing volume of molecular data. Statistical tests for neutrality have been generated exploiting the different properties of the stationary neutral model. Examples of tests are the HKA [2], which takes advantage of the polymorphism/divergence relationship across independent loci in a multilocus framework, and the Lewontin-Krakauer test [3], which looks for an unexpected level of population differentiation in a locus in relation to other loci. Also, there is another family of tests related to linkage disequilibrium, as the one developed by [4], which detect long haplotypes at unusual elevated frequencies in candidate regions.

An important family of these tests, often used as summary statistics, is built on the frequency spectrum of allele polymorphisms. This family includes the well known tests by Tajima [5], Fu and Li [6] and Fay and Wu [7]. If an outgroup is available, these tests are based on the unfolded spectrum ξ_i , that is, the number of segregating sites with a derived allele frequency of i in a sample of (haploid) size n . Without an outgroup, it is not possible to distinguish derived and ancestral alleles and the only available data correspond to the folded spectrum η_i , that is, the number of segregating sites with a minor allele frequency of i . The quantities ξ_i and η_i are all positive and the range of allele frequencies is $1 \leq i \leq n-1$ for the unfolded spectrum, $1 \leq i \leq \lfloor n/2 \rfloor$ for the folded spectrum. The average spectra for the standard Wright-Fisher neutral model are given by

$$E(\xi_i) = \frac{1}{i} \theta L \quad , \quad E(\eta_i) = \frac{n}{i(n-i)(1 + \delta_{i,n-i})} \theta L \quad , \quad (1)$$

where L is the length of the sequence and $\theta = 2p\mu N_e$, where μ is the mutation rate, p is the ploidy and N_e is the effective population size¹. Note that the spectra are proportional to θ .

In a recent paper by Achaz [8], a general framework for these tests was presented. The general tests proposed there were of the form

$$T_\Omega = \frac{\sum_{i=1}^{n-1} i \Omega_i \xi_i}{\sqrt{\text{Var} \left(\sum_{j=1}^{n-1} j \Omega_j \xi_j \right)}} \quad , \quad T_\Omega^* = \frac{\sum_{i=1}^{\lfloor n/2 \rfloor} i \Omega_i^* \eta_i}{\sqrt{\text{Var} \left(\sum_{j=1}^{\lfloor n/2 \rfloor} j \Omega_j^* \eta_j \right)}} \quad (2)$$

¹Note that we define θ as the rescaled mutation rate per base and not per sequence. Apart from this, we follow the notation of [8] and [9].

that are centered (i.e., they have a null expectation value) if the weights Ω_i, Ω_i^* satisfy the conditions $\sum_{i=1}^{n-1} \Omega_i = 0$ and $\sum_{i=1}^{\lfloor n/2 \rfloor} \Omega_i^* = 0$. This framework allows the construction of many new neutrality tests and has been used to obtain optimal tests for specific alternative scenarios [10]. However the original framework does not take into account the dependence of the tests on the sample size, as emphasized in [10]. It is important to choose this dependence in order to obtain results that are as independent as possible on sample size. Moreover, the framework presented in [8] covers just a large subfamily of neutrality tests based on the frequency spectrum, that is, the class of tests that are linear functions of the spectrum. This subfamily contains almost all the tests that can be found in the literature with the exception of the G_ξ, G_η tests of Fu [11], which are second order polynomials in the spectrum whose form is related with Hotelling statistics. Since these G_ξ, G_η tests were shown to be quite effective in some scenarios, it would be interesting to build a general framework for these quadratic (and more generally nonlinear) tests.

In this paper we provide a detailed study of the properties of the whole family of tests based on allele frequency spectrum, beginning with the discussion of the most interesting case, i.e., linear tests. We present a thorough analysis of a simple proposal for the scaling of the tests with the sample size, then we analyze the geometrical properties of the optimal tests presented in [10] and we propose generalizations of D' test to general linear tests and linear optimal tests. Finally, we go beyond the framework presented in [8] and discuss the most general class of tests, that is, polynomials of any order in ξ_i, η_i , and obtain the optimal tests for polynomials of any order. These results allow to build new and more effective tests. The proofs can be found in Appendix B.

2. Linear neutrality tests

2.1. General framework

As discussed by Achaz [8], the general form for linear tests based on the unfolded spectrum can be written as

$$T_\Omega = \frac{\sum_{i=1}^{n-1} i \Omega_i \xi_i}{\sqrt{\text{Var} \left(\sum_{j=1}^{n-1} j \Omega_j \xi_j \right)}} \quad (3)$$

where Ω_i is a set of weights satisfying the condition

$$\sum_{i=1}^{n-1} \Omega_i = 0 . \quad (4)$$

This is the most general form if we require that the test is centered and with variance 1, that is, $E(T_\Omega) = 0$ and $\text{Var}(T_\Omega) = 1$. The condition of centeredness can be obtained substituting the spectrum with its average in the standard neutral model, given by the equations (1).

Alternatively, it is sufficient to choose any pair of unbiased estimators of θ based on the unfolded spectrum

$$\hat{\theta}_\omega = \frac{1}{L} \sum_{i=1}^{n-1} i \omega_i \xi_i \quad , \quad \hat{\theta}_{\omega'} = \frac{1}{L} \sum_{i=1}^{n-1} i \omega'_i \xi_i \quad (5)$$

with weights ω_i, ω'_i that obey the conditions

$$\sum_{i=1}^{n-1} \omega_i = 1 \quad , \quad \sum_{i=1}^{n-1} \omega'_i = 1 \quad (6)$$

to obtain a new test for neutrality:

$$T_\Omega = \frac{\hat{\theta}_\omega - \hat{\theta}_{\omega'}}{\sqrt{\text{Var}(\hat{\theta}_\omega - \hat{\theta}_{\omega'})}} = \frac{\sum_{i=1}^{n-1} i(\omega_i - \omega'_i) \xi_i}{\sqrt{\text{Var}\left(\sum_{j=1}^{n-1} j(\omega_j - \omega'_j) \xi_j\right)}} = \frac{\sum_{i=1}^{n-1} i \Omega_i \xi_i}{\sqrt{\text{Var}\left(\sum_{j=1}^{n-1} j \Omega_j \xi_j\right)}} \quad (7)$$

that is equivalent to the definition (3) with $\Omega_i = \omega_i - \omega'_i$. Therefore a test T_Ω is defined by real vectors Ω or ω, ω' satisfying the above normalization conditions.

If an outgroup is not available, then the test should be based on the folded spectrum and has the general form:

$$\begin{aligned} T_\Omega^* &= \frac{\hat{\theta}_\omega^* - \hat{\theta}_{\omega'}^*}{\sqrt{\text{Var}(\hat{\theta}_\omega^* - \hat{\theta}_{\omega'}^*)}} = \frac{\sum_{i=1}^{\lfloor n/2 \rfloor} i(n-i)(1 + \delta_{n,2i})(\omega_i^* - \omega_{i'}^*) \eta_i}{\sqrt{\text{Var}\left(\sum_{j=1}^{\lfloor n/2 \rfloor} j(n-j)(1 + \delta_{n,2j})(\omega_j^* - \omega_{j'}^*) \eta_j\right)}} = \\ &= \frac{\sum_{i=1}^{\lfloor n/2 \rfloor} i(n-i)(1 + \delta_{n,2i}) \Omega_i^* \eta_i}{\sqrt{\text{Var}\left(\sum_{j=1}^{\lfloor n/2 \rfloor} j(n-j)(1 + \delta_{n,2j}) \Omega_j^* \eta_j\right)}} \quad (8) \end{aligned}$$

where the weights $\Omega_i^* = \omega_i^* - \omega_i^{*'}$ satisfy the conditions

$$\sum_{i=1}^{\lfloor n/2 \rfloor} \omega_i^* = 1 \quad , \quad \sum_{i=1}^{\lfloor n/2 \rfloor} \omega_i^{*'} = 1 \quad \Rightarrow \quad \sum_{i=1}^{\lfloor n/2 \rfloor} \Omega_i^* = 0 \quad (9)$$

and

$$\hat{\theta}_\omega^* = \frac{1}{L} \sum_{i=1}^{\lfloor n/2 \rfloor} \frac{i(n-i)(1+\delta_{n,2i})}{n} \omega_i^* \eta_i \quad , \quad \hat{\theta}_{\omega'}^* = \frac{1}{L} \sum_{i=1}^{\lfloor n/2 \rfloor} \frac{i(n-i)(1+\delta_{n,2i})}{n} \omega_i^{*'} \eta_i \quad (10)$$

are unbiased estimators of θ .

2.2. Transformations of weights and invariance of tests

We report some theorems on the invariance of the tests under affine transformations. These results can be easily proved and are implicitly used throughout this paper.

Theorem 1. *A test of the form (3) does not change its value if all the weights Ω_i are rescaled by a common factor $\lambda > 0$, that is,*

$$\Omega_i \longrightarrow \lambda \Omega_i \quad \Rightarrow \quad T_\Omega \longrightarrow \text{sign}(\lambda) T_\Omega \quad (11)$$

Note that the invariance of the tests mean that these transformations define equivalence classes of weights, i.e., sets of different weights that actually correspond to the same test. In particular, this theorem implies that the space of possible tests, in terms of the weights Ω_i , is not homeomorphic to \mathbb{R}^{n-2} (which would be the subspace of weights in \mathbb{R}^{n-1} that satisfy the linear condition (4)) but to its quotient with respect to the invariance (multiplicative) group \mathbb{R}^+ , that is, the $(n-3)$ -dimensional sphere $S^{n-3} = \mathbb{R}^{n-2}/\mathbb{R}^+$.

Theorem 2. *A test of the form (7) does not change its value under an affine transformation of parameters (λ, ρ_i) on the weights ω_i, ω_i' with a common rescaling factor $\lambda > 0$, that is,*

$$\omega_i \longrightarrow \lambda \omega_i + \rho_i \quad , \quad \omega_i' \longrightarrow \lambda \omega_i' + \rho_i \quad \Rightarrow \quad T_\Omega \longrightarrow \text{sign}(\lambda) T_\Omega \quad (12)$$

However, the estimators (5) are unbiased only if the rescaling factor satisfies the condition $\lambda = 1 - \sum_{i=1}^{n-1} \rho_i$.

2.3. Generalized D' tests for multilocus analysis

The statistic D' [12], which is defined as the ratio of Tajima's D versus its minimum value (given a fixed number of segregating sites), has been used in the literature for multilocus analyses [12, 13, 14], arguing that the value of Tajima's D is affected by the length, the sample size and the number of segregating sites of each studied locus and therefore the values of each locus are not directly comparable.

The contribution of each locus to the heterogeneity is hardly known. Tajima's D is robust to differences in the level of variability (the variance is approximately equal to one) and also quite robust against differences in sample size (as is will be shown in the next part), although the quantitative values of Tajima's D for each condition are somewhat different and therefore the comparison between values is not simple. The proposal of Schaeffer is to use the test

$$D' = \frac{D}{\min(D)_{S=S_{obs}}} \quad (13)$$

as a (re)normalized version of Tajima's D . S_{obs} is the observed number of segregating sites in the sample. This proposal can be generalized for all the tests of the form (3) as follows:

$$T'_\Omega = \frac{T_\Omega}{\min(T_\Omega)_{S=S_{obs}}} = \frac{\sum_{i=1}^{n-1} i\Omega_i\xi_i}{\min(j\Omega_j)S_{obs}} \quad (14)$$

This appears to be the natural generalization of D' to general linear tests.

3. Sample size independent tests

3.1. Scaling of weights with sample size

In this section we would like to remark that there are conditions that have to be imposed on the weights Ω_i or ω_i, ω'_i to ensure that these tests are consistent and meaningful. In fact, the values (and even the number!) of these weights depend explicitly on sample size n . Since every conceivable test should be applied to samples of different size, then its definition involves a whole family of weights $\{\Omega_i^{(n)}\}$ or $\{\omega_i^{(n)}, \omega_i'^{(n)}\}$ with $n = 2, 3 \dots \infty$ and to define a test it is necessary to specify how these weights scale with n .

As an example of the weird effects of some choices of scaling, we consider the test for admixture of [8]. The weights of this test are $\omega_i = \binom{n}{i} 2^{-n} (1 - 2^{-n+1})^{-1}$ and $\omega'_i = 1/(n-1)$. Suppose that the population under study shows

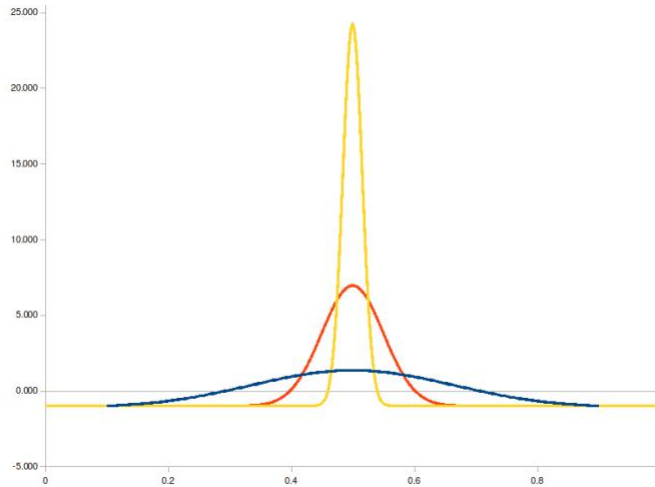


Figure 1: Illustrative example of the dependence of the weight on sample size: weight Ω as a function of i/n for the test for admixture by Achaz, plotted for different sample size $n = 10$ (blue), 100 (red), 1000 (yellow).

an excess of alleles of frequency f between 0.3 and 0.4. The average weight of these frequencies, rescaled by the sample size, is 0.5 for $n = 10$, but it reduces to -0.75 for $n = 100$ and to -1.0 for $n = 1000$. These weights are largely different, even in sign, therefore a strong excess of alleles in this range of frequency would show itself as either a positive or a negative value for this test, depending on the sample size! The reason can be understood by noticing that for n large, the binomial can be approximated by a Gaussian function of the allele frequencies $f = i/n$ centered in $f = 1/2$ and with variance $1/4n$. Therefore this weight function has a strong dependence on n when considered as a function of f and n . The changes of this weight function with sample size are apparent in the plot of Figure 1, which shows the actual function (rescaled by sample size) for $n = 10, 100, 1000$. In this example it is apparent that the interpretation of the results of this test depends on n . This means that the calibration of the test should be different for each possible sample size.

The consistency requirement that we propose is that the result of the test should be almost independent on sample size. This requirement is equivalent to a condition on the scaling of the weights $\Omega_i^{(n)}$ with n . Our proposal for a reasonable requirement on this scaling is the following: the relative weight of

different frequencies in the population should remain approximately constant while varying the size of the sample. This condition ensures that at least for sufficiently large n , the average values of the test on samples of different size from the same population should be approximately independent on sample size, i.e. that the test should be consistent.

To determine the scaling, we note that in limit of large n , the frequency spectrum approaches a continuum and we can define the weights as functions $\Omega(f)$ or $\omega(f), \omega'(f)$ with $f \in (0, 1)$ and $\int_0^1 df \Omega(f) = 0$, $\int_0^1 df \omega(f) = \int_0^1 df \omega'(f) = 1$. Since the ratio of the derived allele count and the sample size i/n is an unbiased estimator of the frequency f of the allele in the population (because $E(i) = nf$), a simple scaling that satisfies the above requirement is

$$\Omega_i^{(n)} \simeq \Omega(i/n) \quad \text{or} \quad \omega_i^{(n)} \simeq \omega(i/n), \quad \omega'_i{}^{(n)} \simeq \omega'(i/n) \quad (15)$$

as proposed by some of the authors in [10].

In order to have the above approximate scaling while obeying the condition $\sum_{i=1}^{n-1} \Omega_i = 0$, there are two simple consistent forms for the weights:

$$\Omega_i^{(n)} = \Omega\left(\frac{i}{n}\right) - \frac{1}{n-1} \sum_{j=1}^{n-1} \Omega\left(\frac{j}{n}\right) \quad (16)$$

where the last term is a (typically small) correction that enforce centeredness of the test, or

$$\Omega_i^{(n)} = \omega_i^{(n)} - \omega'_i{}^{(n)} = \frac{\omega\left(\frac{i}{n}\right)}{\sum_{j=1}^{n-1} \omega\left(\frac{j}{n}\right)} - \frac{\omega'\left(\frac{i}{n}\right)}{\sum_{j=1}^{n-1} \omega'\left(\frac{j}{n}\right)} \quad (17)$$

where the denominators are normalization factors.

Typically this second form (17) for the scaling is more consistent in practice and it is implicitly assumed for most of the existing tests. However, the above expressions give similar numerical results for most choices of the functions $\Omega(f) = \omega(f) - \omega'(f)$. In fact, if $\Omega(f)$ is a limited and piecewise-continuous function, the difference between (16) and (17) is of order $O(\Omega)/n$ (since it is a factor coming from the discretization of the frequencies) and it does not have a relevant impact on the results of the test. Therefore in these cases the two scaling relations (16) and (17) are practically equivalent.

Note that all the tests involving the Watterson estimator (that corresponds to $\omega(f) \sim 1/f$) have additional subtleties that are discussed in the next section.

Example: Fay and Wu's H test

This test was proposed in [7] to look for an excess of high-frequency derived alleles as a signal of selection. It can be defined by the weight functions $\omega(f) = 2f$ and $\omega'(f) = 1$. The weights can be found following equation (17). The resulting test is

$$T_H = \frac{\frac{2}{n(n-1)} \sum_{i=1}^{n-1} i^2 \xi_i - \frac{1}{n-1} \sum_{i=1}^{n-1} i \xi_i}{\sqrt{\text{Var} \left(\frac{2}{n(n-1)} \sum_{j=1}^{n-1} j^2 \xi_j - \frac{1}{n-1} \sum_{j=1}^{n-1} j \xi_j \right)}} \quad (18)$$

The scaling defined in equation (16), with weight function $\Omega(f) = \omega(f) - \omega'(f) = 2f - 1$, gives precisely the same result.

Example: $F(r, r')$ tests of Fu [15]

This large class of test is based on the comparison of two estimators with weights

$$\omega_i = \frac{i^{-r}}{\sum_{j=1}^{n-1} j^{-r}} \quad , \quad \omega'_i = \frac{i^{-r'}}{\sum_{j=1}^{n-1} j^{-r'}} \quad (19)$$

that in the case $r, r' < 1$ correspond precisely to the scaling (17) suggested above, with weight functions $\omega(f) = (1-r)f^{-r}$ and $\omega'(f) = (1-r')f^{-r'}$. This can be easily verified by multiplying both the numerator and the denominator of ω_i, ω'_i by a factor $(1-r)/n^{-r}, (1-r')/n^{-r'}$ respectively. The test by Fay and Wu corresponds actually to $F(-1, 0)$.

The cases with $r \geq 1$ or $r' \geq 1$ involve weight functions with divergent integrals and will be discussed in the next section.

Note that the same weight functions with the scaling (16) would give rise to a slightly different test with weights

$$\Omega_i = (1-r) \left(\frac{i}{n} \right)^{-r} - (1-r') \left(\frac{i}{n} \right)^{-r'} - \left(\frac{(1-r) \sum_{j=1}^{n-1} j^{-r}}{(n-1)n^{-r}} - \frac{(1-r') \sum_{j=1}^{n-1} j^{-r'}}{(n-1)n^{-r'}} \right) \quad (20)$$

that is not consistent for weights of low frequency alleles, i.e. with $i/n \lesssim n^{2/\max(r, r')}$, and therefore less interesting.

Example: test for bottleneck of Achaz [8]

This test is another example of a test with an unwanted scaling:

$$\omega_i = \frac{e^{-\alpha i}}{\sum_{j=1}^{n-1} e^{-\alpha j}} \quad , \quad \omega'_i = \frac{1}{n-1} \quad (21)$$

The weight function for this test is $e^{-\alpha n f} \alpha n / (1 - e^{-\alpha n}) - 1$ that depends strongly on n , therefore this test is not consistent in the above sense.

It is easy to build an equivalent test with the correct scaling by choosing the functions $\omega(f) = \beta e^{-\beta f} / (1 - e^{-\beta})$, $\omega'(f) = 1$. The resulting weights with the scaling (17) are

$$\omega_i = \frac{e^{-\beta i/n}}{\sum_{j=1}^{n-1} e^{-\beta j/n}} = \frac{1 - e^{-\beta/n}}{1 - e^{-\beta(1-1/n)}} e^{-\beta(i-1)/n} \quad , \quad \omega'_i = \frac{1}{n-1} \quad (22)$$

as discussed before. The optimal value reported in [8] is $\alpha \simeq 0.9$ for $n = 30$. This value corresponds to $\beta \simeq 27$.

The test can also be implemented by choosing the scaling (16) and the weight function $\Omega(f) = \omega(f) - \omega'(f) = \beta e^{-\beta f} / (1 - e^{-\beta}) - 1$. The resulting weights are

$$\begin{aligned} \Omega_i &= \frac{\beta e^{-\beta i/n}}{1 - e^{-\beta}} - 1 - \frac{1}{n-1} \left(\frac{\beta(1 - e^{-\beta(1-1/n)})}{(1 - e^{-\beta}) e^{\beta/n} (1 - e^{-\beta/n})} - (n-1) \right) = \\ &= \frac{\beta(1 - e^{-\beta(1-1/n)})}{(1 - e^{-\beta}) e^{\beta/n} (1 - e^{-\beta/n})} \cdot \left(\frac{1 - e^{-\beta/n}}{1 - e^{-\beta(1-1/n)}} e^{-\beta(i-1)/n} - \frac{1}{n-1} \right) \end{aligned} \quad (23)$$

that are equivalent to the weights (22) up to an irrelevant multiplicative factor (see Theorem 1). Therefore in this case the two choices of scaling give precisely the same result.

3.2. Divergent weights

As discussed above, the two choices of scaling in equation (16) and (17) do not usually bring to sensibly different numerical results. However, there are important choices of $\Omega(f)$ for which this approximate equivalence between (16) and (17) does not hold. These critical cases correspond to functions that diverge as $1/f$ or faster near $f = 0$ (or $f = 1$). This divergence is not a real feature of the distribution, because the integral has a natural cutoff at the scale of the inverse population size² $f_{min} = 1/N$, but in this case the integral $\int_{1/N}^1 df \Omega(f)$ has a strong dependence on the cutoff $1/N$ and therefore the function $\Omega(f)$ itself should depend strongly on N to ensure proper normalization.

²Or more precisely the effective population size $1/N_e$, but this does not affect the above discussion.

If this dependence is contained in an multiplicative term in front of $\omega(f)$ or $\omega'(f)$ or both, then the second term in equation (16) is not a small correction of order $1/n$ as it happens with simple functions $\Omega(f)$, but rather it represents a relevant correction with a strong dependence on sample size n and population size N . The denominators in equation (17) also show a strong dependence on n (that could not be avoided anyway) but not on N , and therefore this second scaling form should be used. The dependence on sample size is as strong as the dependence of the divergent integral from the cutoff³: for functions diverging as f^{-k} with $k \geq 1$, the dependence on n goes as n^{1-k} if $k > 1$ or $\log(n)$ for $k = 1$. This case always occurs when the test is build by comparing an estimator of θ with the Watterson estimator, which corresponds to $\omega(f) \sim 1/f$ and therefore has a logarithmic dependence on n given by the usual harmonic factor $a_n = \sum_{j=1}^{n-1} 1/j \simeq \log(n) + \gamma + O(1/n)$. A well-known examples of this case is Tajima's D [5].

If the dependence of $\Omega(f)$ on N is contained in an additive term that does not depend on f , it is the correction in (16) that does not depend on N and therefore the first scaling form is more appropriate. We do not know examples of tests of this kind in the literature, even if the test by Zeng *et al.* [16] can be interpreted also in this way.

Example: Tajima's D test

This is the most known test for neutrality based on the frequency spectrum. It is given by the difference between the Tajima estimator Π [17] based on the nucleotide pairwise diversity Π and the Watterson estimator θ_W [18] based on the number S of segregating sites, therefore it can be defined by the weight functions $\omega(f) = 2(1 - f)$ for Π and $\omega'(f) = 1/f \log(N)$ for the Watterson estimator. The latter function has an integral that diverges logarithmically near $f = 0$, and the corresponding dependence on N is contained in the factor $1/\log(N)$ that multiplies $\omega'(f)$, therefore the scaling (17) should

³This can be easily understood by noticing that the sample size n plays the role of the cutoff in the sum over the frequencies that are present in the sample, which is the same role played by the population size N for the whole population. More formally, the denominator in equation (17) can be bounded from above and from below by the divergent integral, and therefore the divergence of the denominator as $n \rightarrow \infty$ will be the same as the divergence of the integral as its inverse cutoff (that is, N) goes to infinity.

be used. The result is the usual test

$$T_D = \frac{\sum_{i=1}^{n-1} \frac{2i(n-i)}{n(n-1)} \xi_i - S/a_n}{\sqrt{\text{Var} \left(\sum_{j=1}^{n-1} \frac{2j(n-j)}{n(n-1)} \xi_j - S/a_n \right)}} = \frac{\Pi - S/a_n}{\sqrt{\text{Var} (\Pi - S/a_n)}} \quad (24)$$

Example: test of Zeng et al. [16]

This test was proposed to look for an excess of high-frequency derived alleles compared to low-frequency alleles. It is defined by the weight functions $\omega(f) = 1$ and $\omega'(f) = 1/f \log(N)$, the latter corresponding to the Watterson estimator. Proceeding as in the above example, the result is

$$T_E = \frac{\sum_{i=1}^{n-1} \frac{i}{(n-1)} \xi_i - S/a_n}{\sqrt{\text{Var} \left(\sum_{j=1}^{n-1} \frac{j}{(n-1)} \xi_j - S/a_n \right)}} \quad (25)$$

Note that exceptionally the scaling of this test can also be defined by (16), without modifying the result. This is a consequence of the two equivalent forms for the weight function, $\Omega(f) = 1 - 1/f \log(N)$ or $\Omega(f) = \log(N) - 1/f$.

3.3. Weights of singletons

The above scaling (15) is valid in principle for all weights. However in practice there is an important exception, that is, the weight Ω_1 of singletons. This is due to the fact that for $n \ll N$, the number of derived singletons ξ_1 is the only estimator that is affected by very rare derived alleles (and often by sequencing errors, see [19]). More precisely, ξ_1 is actually the only estimator sensitive to the deviations from neutrality in alleles of frequency $1/N < f < 1/n$, which represent a vast majority of the SNPs in the population and can contain interesting biological information. Therefore, if the contribution of these alleles is relevant for the test, we can enhance (or reduce) the weight Ω_1 by adding a factor Ω_{ds} .

In the approach detailed in the previous sections, this additional contribution to Ω_1 is needed to take into account a contribution $\Delta\Omega(f)$ to $\Omega(f)$ of the form $\Delta\Omega(f) = \Omega_{ds}I(f < \phi)/\phi$ with $\phi \ll 1$. As far as the maximum sample size never exceeds in practice $n_{max} \ll 2/\phi$, this function weights positively only alleles that appear as singletons.

Similarly, ω_1 and ω'_1 can be enhanced by ω_{ds} , ω'_{ds} that correspond to contributions $\Delta\omega(f) = \omega_{ds}I(f < \phi)/\phi$, $\Delta\omega'(f) = \omega'_{ds}I(f < \phi)/\phi$. The test of Fu and Li [6] fall into this case.

A similar argument applies also to the weights of the number of ancestral singletons, that is, Ω_{n-1} , ω_{n-1} , ω'_{n-1} that can be enhanced by factors Ω_{as} , ω_{as} and ω'_{as} respectively. However this case is more rare, the only interesting example being the tests of Achaz [19] that avoid sequencing errors by neglecting both derived and ancestral singletons.

Summarizing the results up to this section, a test T_Ω is completely defined by a function $\Omega(f)$ and two parameters Ω_{ds} , Ω_{as} (that could depend on n) satisfying the conditions

$$\Omega_{ds} + \Omega_{as} + \int_0^1 df \Omega(f) = 0 \quad (26)$$

and determining the weights through the formula:

$$\Omega_i^{(n)} = \Omega\left(\frac{i}{n}\right) + \Omega_{ds}\delta_{i,1} + \Omega_{as}\delta_{i,n-1} - \frac{1}{n-1} \left(\Omega_{ds} + \Omega_{as} + \sum_{j=1}^{n-1} \Omega\left(\frac{j}{n}\right) \right) \quad (27)$$

or by a pair of functions $\omega(f)$, $\omega'(f)$ and parameters ω_{ds} , ω'_{ds} , ω_{as} , ω'_{as} satisfying

$$\omega_{ds} + \omega_{as} + \int_0^1 df \omega(f) = \omega'_{ds} + \omega'_{as} + \int_0^1 df \omega'(f) = 1 \quad (28)$$

and resulting in this formula for the scaling of the weights:

$$\Omega_i^{(n)} = \frac{\omega_{ds}\delta_{i,1} + \omega_{as}\delta_{i,n-1} + \omega\left(\frac{i}{n}\right)}{\omega_{ds} + \omega_{as} + \sum_{j=1}^{n-1} \omega\left(\frac{j}{n}\right)} - \frac{\omega'_{ds}\delta_{i,1} + \omega'_{as}\delta_{i,n-1} + \omega'\left(\frac{i}{n}\right)}{\omega'_{ds} + \omega'_{as} + \sum_{j=1}^{n-1} \omega'\left(\frac{j}{n}\right)} \quad (29)$$

As showed in the examples above and below, most of the tests in the literature have this general scaling, with the only exceptions of the ones contained in [8] that are not consistent in the above sense.

Example: Fu and Li's F test

This test looks for an excess of very rare derived alleles as a possible signature of negative selection [6]. The only nonzero weights are $\omega_{ds} = 1$ and $\omega'(f) = 1/f \log(N)$, while $\omega(f) = \omega'_{ds} = \omega_{as} = \omega'_{as} = 0$. The resulting test is

$$T_F = \frac{\xi_1 - S/a_n}{\sqrt{\text{Var}(\xi_1 - S/a_n)}} \quad (30)$$

Note that this test has both singleton weights and a divergent weight function.

Example: error-corrected tests of Achaz [19]

This class of tests is an attempt to correct for sequencing errors and biases in the data by removing the alleles where most of the problems manifest themselves, i.e. singletons (both ancestral and derived). With a slight generalization of the proposal in [19], the weights of the singletons are chosen in such a way to cancel precisely the contributions of the weight functions:

$$\Omega_{ds} = -\Omega\left(\frac{1}{n}\right), \Omega_{as} = -\Omega\left(1 - \frac{1}{n}\right) \quad (31)$$

or

$$\omega_{ds} = -\omega\left(\frac{1}{n}\right), \omega_{as} = -\omega\left(1 - \frac{1}{n}\right), \omega'_{ds} = -\omega'\left(\frac{1}{n}\right), \omega'_{as} = -\omega'\left(1 - \frac{1}{n}\right) \quad (32)$$

therefore the final weights of derived or ancestral singletons are zero. These corrections can be applied in principle to any weight function.

3.4. Scaling of weights in tests without an outgroup

The above arguments can be repeated in a straightforward way for the tests T_{Ω}^* based on the folded spectrum η_i . The only relevant difference is that the frequency f of the minor allele in the population is always less than 50%, that is, $f \in (0, 1/2]$. For consistency with the unfolded case, the weight $\eta_{n/2}$ is reduced by a factor 2. Moreover, the additional parameters related to the weights of singletons cannot distinguish between ancestral and derived alleles and therefore reduce to Ω_s^* , ω_s^* , $\omega_s^{*'}$. These parameter, together with the functions $\Omega^*(f)$, $\omega^*(f)$ and $\omega^{*'}(f)$, should satisfy the conditions

$$\Omega_s^* + \int_0^{1/2} df \Omega^*(f) = 0 \quad , \quad \omega_s^* + \int_0^{1/2} df \omega^*(f) = \omega_s^{*'} + \int_0^{1/2} df \omega^{*'}(f) = 1 \quad (33)$$

The formulae that determine the scaling of the weights are:

$$\Omega_i^{*(n)} = \frac{1}{1 + \delta_{n,2i}} \Omega^*\left(\frac{i}{n}\right) + \Omega_s^* \delta_{i,1} - \frac{1}{\lfloor n/2 \rfloor} \left(\Omega_s^* + \sum_{j=1}^{\lfloor n/2 \rfloor} \frac{1}{1 + \delta_{n,2j}} \Omega^*\left(\frac{j}{n}\right) \right) \quad (34)$$

$$\Omega_i^{*(n)} = \frac{\omega_s^* \delta_{i,1} + \omega^*\left(\frac{i}{n}\right) / (1 + \delta_{n,2i})}{\omega_s^* + \sum_{j=1}^{\lfloor n/2 \rfloor} \omega^*\left(\frac{j}{n}\right) / (1 + \delta_{n,2j})} - \frac{\omega_s^{*'} \delta_{i,1} + \omega^{*'}\left(\frac{i}{n}\right) / (1 + \delta_{n,2i})}{\omega_s^{*'} + \sum_{j=1}^{\lfloor n/2 \rfloor} \omega^{*'}\left(\frac{j}{n}\right) / (1 + \delta_{n,2j})} \quad (35)$$

The weights of the folded versions of Tajima's D and Fu and Li's F^* test follow this scaling. The nonzero weight functions are $\omega^*(f) = 1$, $\omega'^*(f) = 1/(\log(N)f(1-f))$ for Tajima's D and $\omega_s^* = 1$, $\omega'^s(f) = 1/(\log(N)f(1-f))$ for the test of Fu and Li.

3.5. Alternative choices of scaling

The choice of scaling discussed in the previous sections represents a quite simple and effective way to fix the dependence on n of a newly devised test. However, other choices are possible whose weights differ from the above ones for small n . The reason is that for n not too large, both the variance of order $f(1-f)/n \simeq i(n-i)/n^3$ in the estimation of the frequency $f = i/n$ and the related uncertainty about how the frequencies are actually weighted in the test become important. This uncertainty originates from the (binomial) sampling of individuals from the population and there is some degree of arbitrariness in deciding how to account for it. Moreover, tests that take it into account could be not consistent in the above sense.

A possible choice of scaling that uses the binomial sampling is the following: considering $\omega(f)$, $\omega'(f)$ as frequency distributions, the weights ω_i , ω'_i are assigned from $\omega(f)$, $\omega'(f)$ through the same binomial sampling that is done for allele spectra, that is,

$$\omega_i = \frac{\int_0^1 df \binom{n}{i} f^i (1-f)^{n-i} \omega(f)}{\int_0^1 df (1-f^n - (1-f)^n) \omega(f)} \quad (36)$$

$$\omega'_i = \frac{\int_0^1 df \binom{n}{i} f^i (1-f)^{n-i} \omega'(f)}{\int_0^1 df (1-f^n - (1-f)^n) \omega'(f)} \quad (37)$$

A simple example of this scaling (but with an highly divergent weight function) is given by the test for admixture [8] discussed before. Optimal tests also follow this scaling.

Example: test for admixture of Achaz [8]

This test is apparently not consistent and it does not follow the scaling (15). However it follows another scaling related to the allele sampling. To understand this, consider the weight functions $\omega(f) = \delta(f - 1/2)$, $\omega'(f) = 1$ where $\delta(f - 1/2)$ is a Dirac delta function⁴ centered in $1/2$. If we scale the

⁴The Dirac delta $\delta(f - a)$ is a function whose value is 0 if $f \neq a$ and $+\infty$ if $f = a$. The integral $\int \delta(f - a)g(f)df$ is $g(a)$ if a is inside the range of integration and 0 otherwise.

weights according to (36),(37), that is,

$$\omega_i = \frac{\int_0^1 df \binom{n}{i} f^i (1-f)^{n-i} \omega(f)}{\int_0^1 df (1-f^n - (1-f)^n) \omega(f)} = \frac{\binom{n}{i} 2^{-n}}{1 - 2^{-n+1}} \quad (38)$$

$$\omega'_i = \frac{\int_0^1 df \binom{n}{i} f^i (1-f)^{n-i} \omega'(f)}{\int_0^1 df (1-f^n - (1-f)^n) \omega'(f)} = \frac{1}{n-1} \quad (39)$$

then the corresponding test is precisely the one proposed by Achaz. Note that the strong dependence of the test from sample size does not come only from the choice of scaling, but also from the weight function chosen, that is highly divergent.

4. Linear optimal tests

4.1. On the existence of generic tests

An interesting question on the way to build good linear tests is the following: do there exist generic tests? A completely generic test for neutrality should be able to detect any deviation from the spectrum of the null model that is sufficiently large. Unfortunately, these tests do not exist. In fact, for every test defined by a set of weights Ω_i it is possible to find a spectrum $\xi_i = \alpha/ia_n + (1-\alpha)\Delta_i$ that is maximally different from the standard spectrum at least in a range of frequencies and is nevertheless undetectable by the test because its average value on this spectrum is zero. This is expressed in a more formal way in the following theorem, which shows that even the complete lack of alleles in some range of frequencies could not be always detected.

Theorem 3. *For every set of n real weights Ω_i with $\sum_i \Omega_i = 0$, there is a set of n real numbers $\Delta_i \neq \text{const}/i$ and a parameter $\alpha \in [0, 1]$ that satisfy the conditions*

$$\sum_i i\Omega_i \Delta_i = 0 \quad , \quad \min_{i \in [1, n-1]} \left(\alpha \frac{1}{ia_n} + (1-\alpha)\Delta_i \right) = 0 \quad (40)$$

Actually this function is not a mathematical function, but a distribution, i.e. an element of a dual space of regular functions.

The above limitation is not a consequence of the small sample size. This can be seen for example in the framework of the scaling theory discussed in this paper. In fact, for large sample size, the weights can be approximated by a weight function $\Omega(f)$. In this context it is possible to prove the next theorem, that is a continuous equivalent of the previous one.

Theorem 4. *For every piece-wise continuous weight function $\Omega(f) \in L^1_{[1/N,1]}$ such that $\int_{1/N}^1 \Omega(f)df = 0$, there is a smooth function $\Delta(f) \neq \text{const}/f$ and a parameter $\alpha \in [0, 1]$ that satisfy the conditions*

$$\int_{1/N}^1 df f \Omega(f) \Delta(f) = 0 \quad , \quad \inf_{f \in [0,1]} \left(\alpha \frac{1}{f \log(N)} + (1 - \alpha) \Delta(f) \right) = 0 \quad (41)$$

Note that in principle this problem can be solved using multiple tests. In fact multiple tests should be able to detect any strong deviation from the null spectrum, provided that the number of these tests is large enough, as can be seen from the following theorem.

Theorem 5. *Given at least $n - 2$ linearly independent sets of $n - 1$ real weights Ω_i with $\sum_i \Omega_i = 0$, it is not possible to find a set of real numbers $\Delta_i \neq \text{const}/i$ such that $\sum_i i \Omega_i \Delta_i = 0$.*

This last theorem is only a formal result and the requirement of $n - 2$ independent tests is too strong. In practice a small (but good) set of tests can detect most of the reasonable and interesting deviations for realistic spectra.

The above theorems can be extended to the folded spectrum. In this section and the next ones, we will consider only tests based on the unfolded spectrum. The generalization of the discussion to the folded spectrum is usually straightforward after substituting ξ_i ($i = 1 \dots n - 1$) with η_i ($i = 1 \dots \lfloor n/2 \rfloor$).

4.2. Optimal tests and their geometric structure

From the theorems of the previous section, it is clear that a single test cannot detect all the possible deviations occurring in complicated evolutionary scenarios. However it is still possible to optimize neutrality tests of for a specific alternative evolutionary scenario. A simple optimality condition has been proposed by some of the authors in [10] in order to maximize the power of the test to detect a fixed alternative scenario. If the null spectrum is $E(\xi_i) = \theta L \xi_i^0$ and the expected spectrum of the alternative scenario

is $\mathcal{E}(\xi_i) = \theta L \bar{\xi}_i$, the condition for optimal tests is the maximization of the average result of the test under the alternative scenario:

$$\mathcal{E}(T_\Omega) = \frac{\sum_{i=1}^{n-1} \Omega_i \theta L \bar{\xi}_i / \xi_i^0}{\sqrt{\text{Var} \left(\sum_{j=1}^{n-1} \Omega_j \xi_j / \xi_j^0 \right)}} \quad (42)$$

This condition is based on the observation that the tests have mean zero and variance 1, therefore if the distributions of the results of the tests are similar, the maximization of the average value of the test should correspond to the maximization of the average power of the test. It is also possible to maximize directly the power of the test, taking into account the different distribution of the results under the null and the alternative model; this possibility will be pursued in section 4.5.

Interestingly, optimal tests show a geometric structure which becomes apparent after defining the scalar product between spectra:

$$\langle\langle \xi', \xi'' \rangle\rangle \equiv \sum_{i,j} \xi'_i c_{ij}^{-1} \xi''_j \quad (43)$$

where c_{ij}^{-1} is the inverse of the covariance matrix $\text{Cov}(\xi_i, \xi_j)$. Since $\text{Cov}(\xi_i, \xi_j)$ is symmetric and positive, its inverse is also symmetric and positive, i.e. it is a positive bilinear form, therefore the above expression defines a scalar product. Then the optimal test for an alternative spectrum $\bar{\xi}$ can be written in the elegant form⁵

$$T_O = \frac{\langle\langle \xi, \bar{\xi} \rangle\rangle - \langle\langle \xi, \xi^0 \rangle\rangle \langle\langle \xi^0, \bar{\xi} \rangle\rangle / \langle\langle \xi^0, \xi^0 \rangle\rangle}{\sqrt{\langle\langle \bar{\xi}, \bar{\xi} \rangle\rangle - \langle\langle \xi^0, \bar{\xi} \rangle\rangle^2 / \langle\langle \xi^0, \xi^0 \rangle\rangle}} \quad (44)$$

The numerator of the test is actually the matrix element between $\bar{\xi}$ and ξ of the linear operator $1 - P_{\xi^0}$, where P_{ξ^0} is the projection operator along ξ^0 . In other words, it is proportional to the difference between the length of the projection of ξ on $\bar{\xi}$ and the length of the projection on $\bar{\xi}$ of the spectrum obtained by the projection of ξ on ξ^0 , as illustrated in Figure 2.

⁵We do not provide a proof of this expression here because it can be easily obtained as a special case of the general formula (67) that we will discuss later in the context of nonlinear tests. A direct proof of this result can be found in [10] after substituting the scalar products with the definition (43).

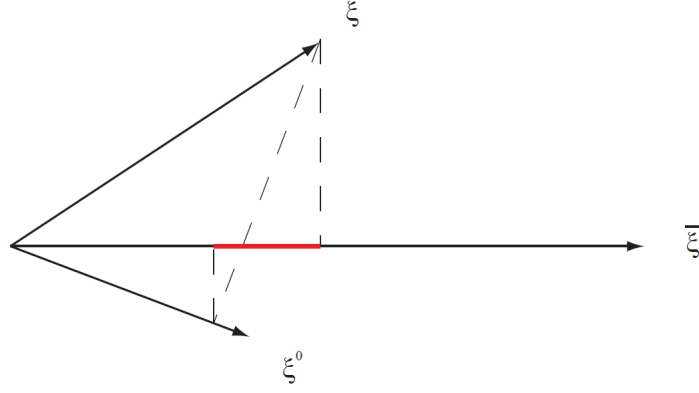


Figure 2: Geometrical representation of the numerator of the optimal test T_O in (44). The length of the red line segment corresponds to the value of the numerator.

From this geometrical interpretation it is clear that if the spectrum ξ corresponds to the null spectrum ξ^0 , then the two projections are equal and the result of the test is zero. On the other side, if the spectrum is the alternative spectrum $\bar{\xi}$, then the value of the test is

$$T_O^{(max)} = \theta L \sqrt{\langle \langle \bar{\xi}, \bar{\xi} \rangle \rangle - \frac{\langle \langle \xi^0, \bar{\xi} \rangle \rangle^2}{\langle \langle \xi^0, \xi^0 \rangle \rangle}} \quad (45)$$

which is the maximum value over all possible tests in the alternative scenario. The same expression, but with a minus sign, corresponds to the minimum value.

The denominator of the test is the square root of the matrix element of the linear operator $1 - P_{\xi^0}$ between $\bar{\xi}$ and itself. Note that both the numerator and the denominator of the test do not change by adding any (possibly negative) multiple of ξ^0 to $\bar{\xi}$, because ξ^0 lies in the kernel of $1 - P_{\xi^0}$. This means that optimal tests depend only on the expected deviations from the null spectrum in the alternative scenario. The result of the test is maximum when the deviations of the data from the null spectrum correspond exactly to the expected ones, and it is minimum when they are opposite to the expected ones.

4.3. Scaling of optimal tests

Optimal tests have weights proportional to the expected allele distribution $\omega_i = \bar{\xi}_i / \sum_{j=1}^{n-1} \bar{\xi}_j$ and to the null allele distribution $\omega'_i = \xi_i^0 / \sum_{j=1}^{n-1} \xi_j^0$.

Therefore the weights of an optimal test follow the same scaling with sample size as the allele distributions. Denoting by $\bar{\xi}(f)$ and $\xi^0(f)$ the spectra of expected and null allele frequencies in the whole population, the spectra for the sample are obtained by binomial sampling

$$\xi_i = \int_{1/N_e}^1 df P_{bin}(i; n, f) \xi(f) \quad P_{bin}(i; n, f) = \binom{n}{i} f^i (1-f)^{n-i} \quad (46)$$

from the spectra $\xi(f) = \bar{\xi}(f)$ and $\xi(f) = \xi^0(f)$ respectively. Therefore the scaling of the allele distributions is

$$\frac{\bar{\xi}_i}{\sum_{j=1}^{n-1} \bar{\xi}_j} = \frac{\int_{1/N_e}^1 df \binom{n}{i} f^i (1-f)^{n-i} \bar{\xi}(f)}{\int_{1/N_e}^1 df (1-f^n - (1-f)^n) \bar{\xi}(f)} \quad (47)$$

$$\frac{\xi_i^0}{\sum_{j=1}^{n-1} \xi_j^0} = \frac{\int_{1/N_e}^1 df \binom{n}{i} f^i (1-f)^{n-i} \xi^0(f)}{\int_{1/N_e}^1 df (1-f^n - (1-f)^n) \xi^0(f)} \quad (48)$$

which is precisely the scaling (36),(37) with weight functions $\omega(f) \propto \bar{\xi}(f)$ and $\omega'(f) \propto \xi^0(f)$. This scaling does not correspond to the scaling (15) suggested in this paper, but it takes into account the sampling process in a straightforward way, being based on the expected and null allele distributions for the sample and therefore immediately related to the binomial sampling of alleles from the population.

Note that these weights actually follow the scaling (15) for large n , with the same weight functions $\omega(f) \propto \bar{\xi}(f)$ and $\omega'(f) \propto \xi^0(f)$. This agrees with the fact that when the sample size is large enough, the variance of the sampling process can be safely ignored and all reasonable choices of scaling are equivalent to our proposal (15).

4.4. Generalizing D' for optimal test

The D' statistics of [12] can be generalized for optimal tests as it was done for general linear tests. In particular, for a fixed number of segregating sites S_{obs} , the generalization for optimal tests in the approximation of unlinked sites and $\theta \ll 1$ is

$$T'_O = \frac{\sum_{i=1}^{n-1} \Omega_i \xi_i / \xi_i^0}{\min_j (\Omega_j / \xi_j^0) S_{obs}} = \frac{\sum_{i=1}^{n-1} (\xi_i / S_{obs}) (\bar{\xi}_i / \sum_{j=1}^{n-1} \bar{\xi}_j) / (\xi_i^0 / \sum_{j=1}^{n-1} \xi_j^0) - 1}{\min_k ((\bar{\xi}_k / \sum_{l=1}^{n-1} \bar{\xi}_l) / (\xi_k^0 / \sum_{l=1}^{n-1} \xi_l^0)) - 1} \quad (49)$$

which has the interesting property of depending only on the allele frequency distributions.

However in the case of optimal tests there is another possibility, namely to define a \bar{D}' test as the ratio of the optimal test and of its average minimum, assuming that the spectrum corresponds to the average spectrum of the actual scenario. This would correspond to the form

$$\bar{T}'_O = - \frac{\langle\langle\xi, \bar{\xi}\rangle\rangle - \langle\langle\xi, \xi^0\rangle\rangle\langle\langle\xi^0, \bar{\xi}\rangle\rangle/\langle\langle\xi^0, \xi^0\rangle\rangle}{S_{obs}/a_n \sqrt{(\langle\langle\bar{\xi}, \bar{\xi}\rangle\rangle - \langle\langle\xi^0, \bar{\xi}\rangle\rangle^2/\langle\langle\xi^0, \xi^0\rangle\rangle) (\langle\langle\xi, \xi\rangle\rangle - \langle\langle\xi^0, \xi\rangle\rangle^2/\langle\langle\xi^0, \xi^0\rangle\rangle)}} \quad (50)$$

which has the interesting property of being symmetric with respect to the actual spectrum ξ and the expected spectrum $\bar{\xi}$.

4.5. Linear tests with maximum power

The condition for optimal tests is the maximization of $\mathcal{E}(T_\Omega)$ under the alternative scenario. However, a better approach would be the maximization of the power of the test to reject the neutral model in the alternative scenario, given a choice of significance level α . This approach requires the knowledge of the form of the probability distributions $p(T_\Omega = t|H_0)$, $p(T_\Omega = t|H_1)$ where H_0 and H_1 are the null and alternative model, or equivalently of all the moments of the spectrum $E(\xi_i \xi_j \xi_k \dots)$ and $\mathcal{E}(\xi_i \xi_j \xi_k \dots)$.

Since this information is usually not available in analytic form and hard to obtain computationally, we limit to the case where the distribution can be well approximated by a Gaussian both for the null and for the expected model. Then the only information needed are the spectra $\mu_i = E(\xi_i)$, $\bar{\mu}_i = \mathcal{E}(\xi_i)$ and their covariance matrices $c_{ij} = E(\xi_i \xi_j) - E(\xi_i)E(\xi_j)$, $\bar{c}_{ij} = \mathcal{E}(\xi_i \xi_j) - \mathcal{E}(\xi_i)\mathcal{E}(\xi_j)$.

We expect that both in this approximation and in the general case, the tests with maximum power will depend on the significance level chosen, therefore limiting the interest of these tests and the possibilities of comparison between results of the test on samples from different experiments.

We call $\tau = \text{erf}^{-1}(1 - 2\alpha)$ the z -value corresponding to the critical p -value α . In the Gaussian approximation, the power is given by the following expression

$$\text{Power} = \frac{1}{2} \left[1 + \text{erf} \left(\frac{\sum_j \bar{\mu}_j \Omega_j - \tau \sqrt{\sum_{j,k} c_{jk} \Omega_j \Omega_k}}{\sqrt{\sum_{j,k} \bar{c}_{jk} \Omega_j \Omega_k}} \right) \right] \quad (51)$$

then its maximization is equivalent to the maximization of

$$\frac{\sum_j \bar{\mu}_j \Omega_j - \tau \sqrt{\sum_{j,k} c_{jk} \Omega_j \Omega_k}}{\sqrt{\sum_{j,k} \bar{c}_{jk} \Omega_j \Omega_k}} \quad (52)$$

In the general case, the weights corresponding to the maximum depend explicitly on τ and therefore on α . This dependence is expected but unwanted, since the interpretation of the test depends explicitly on the critical p -value chosen. Obtaining explicit solutions for the weights is complicated and will not be discussed here.

There is only one case with weights independent on τ , that is the case of \bar{c}_{ij} (approximately) proportional to c_{ij} . In this case the maximization of the power of the test is (approximately) equivalent to the maximization of the average result of the test, which is precisely our condition for optimal tests.

There is also a regime of values of α such that the weights corresponding to maximum power are independent on α , that is, the regime $\tau(\alpha) \gg 1$. In this case the power is an increasing function of $\sum_{j,k} \bar{c}_{jk} \Omega_j \Omega_k / \sum_{l,m} c_{lm} \Omega_l \Omega_m$ and the weights are simply given by the null eigenvector (or linear combination of null eigenvectors) of the matrix $\bar{c}_{ij} - \chi c_{ij}$, where χ is uniquely defined by the requirement that $\bar{c}_{ij} - \chi c_{ij}$ be a negative semidefinite matrix with at least a null eigenvalue. However, this regime is uninteresting because such small significance levels are practically useless (if $\tau \sim 10$, the corresponding critical p -value is $\alpha \sim 10^{-20}$).

In our opinion, the maximization of (52) is not interesting in practice because of the dependence on α and of the complicated form of the corresponding weights. Optimal tests represent a good compromise between high power, simplicity and easiness to interpret and compare the results. However, it could be possible to build interesting tests with higher power by selecting a linear combinations of the weights of the two α -independent tests discussed above, that is, optimal tests and tests that maximize the alternative/null variance ratio $\sum_{j,k} \bar{c}_{jk} \Omega_j \Omega_k / \sum_{l,m} c_{lm} \Omega_l \Omega_m$.

5. Beyond linear neutrality tests

5.1. Quadratic and nonlinear tests

Almost all the neutrality tests proposed in the literature are linear in the spectrum ξ_i . As far as we know, there is only one exception, namely the G_ξ

test of Fu [11]. This test is a quadratic polynomial reminiscent of Hotelling's t^2 statistics for the different components of the spectra:

$$G = \sum_{i,j=1}^{n-1} c_{ij}^{-1} (\xi_i - \theta L \xi_i^0) (\xi_j - \theta L \xi_j^0) \quad (53)$$

where c_{ij}^{-1} is the inverse of the covariance matrix $\text{Cov}(\xi_i, \xi_j)$. Actually the test proposed by Fu is an approximation to this test with a different normalization, namely

$$G_\xi = \frac{1}{n-1} \sum_{i=1}^{n-1} \frac{(\xi_i - \theta L \xi_i^0)^2}{\text{Var}(\xi_i)} \quad (54)$$

In this approximation the off diagonal terms in the covariance can be neglected [9, 11]. For large samples, the distribution of the results of the test G tends to a χ^2 distribution with $n-1$ degrees of freedom.

Fu's approach cannot be extended to general quadratic or higher order tests, because the distribution of the results of the test would be generally unknown and not positive definite. For this reason we propose to rescale the tests to have zero mean and variance 1. With this normalization, we expect that the distribution would asymptotically converge to a Gaussian $N(0, 1)$ for all tests. As an example, the (re)normalized version of Fu's test would be

$$T_G = \frac{\sum_{i,j=1}^{n-1} c_{ij}^{-1} (\xi_i - \theta L \xi_i^0) (\xi_j - \theta L \xi_j^0) - (n-1)}{\sqrt{\text{Var} \left(\sum_{i,j=1}^{n-1} c_{ij}^{-1} (\xi_i - \theta L \xi_i^0) (\xi_j - \theta L \xi_j^0) - (n-1) \right)}} \quad (55)$$

Since the only difference between this test and the original one is the normalization and a shift by a constant factor $n-1$, the power of the test is the same.

Now we present a systematic discussion of nonlinear tests that are generic polynomials (or eventually power series) in the spectrum ξ_i . All the tests are rescaled to be centered (i.e., to have zero mean) and have variance 1. We denote by $\mu_{ijk\dots}$ the moments of the spectrum under the null model, that is, $\mu_{ijk\dots} = E(\xi_i \xi_j \xi_k \dots)$. With this definition, $\mu_i = \theta L \xi_i^0$. Note that all these moments depend on θ . In the approximation of unlinked (independent) sites and small θ , the second moments are equal to $\mu_{ij} = \theta L \xi_i^0 \delta_{ij} + \theta^2 L^2 \xi_i^0 \xi_j^0$.

The weights of general nonlinear tests can depend explicitly on θ , as seen in the previous example. To compute the values of the tests, the (unknown)

parameter θ is substituted with an estimator $\hat{\theta}$. Unlike the linear case, in this case there are two different classes of tests, related to the dependence on $\hat{\theta}$ of the centeredness: *strongly centered* and *weakly centered* tests.

Strongly centered tests are tests that are always centered for any value of $\hat{\theta}$, even if it is different from the actual value of θ . The general form for strongly centered tests is

$$T_{\Omega} = \frac{\sum_{i=1}^{n-1} \Omega_i^{(1)} \xi_i + \sum_{i,j=1}^{n-1} \Omega_{ij}^{(2)} \xi_i \xi_j + \sum_{i,j,k=1}^{n-1} \Omega_{ijk}^{(3)} \xi_i \xi_j \xi_k + \dots}{\sqrt{\text{Var} \left(\sum_{i=1}^{n-1} \Omega_i^{(1)} \xi_i + \sum_{i,j=1}^{n-1} \Omega_{ij}^{(2)} \xi_i \xi_j + \sum_{i,j,k=1}^{n-1} \Omega_{ijk}^{(3)} \xi_i \xi_j \xi_k + \dots \right)}} \quad (56)$$

with the real symmetric weights $\Omega_{ijk\dots}^{(n)}$ satisfying the set of conditions

$$0 = \sum_{i=1}^{n-1} \Omega_i^{(1)} \mu_i^{(m)} + \sum_{i,j=1}^{n-1} \Omega_{ij}^{(2)} \mu_{ij}^{(m)} + \dots \quad , \quad m = 1, 2, 3 \dots \quad (57)$$

where we denote by $\mu_{ijk\dots}^{(p)}$ the p -th term of the Taylor expansion with respect to θL of $\mu_{ijk\dots}$ ⁶. The sum can be limited to polynomials of some finite order in ξ_i or it can be a (convergent) power series. If we introduce the notation $\mathbf{I} = ijk\dots$ to denote a group of $n_{\mathbf{I}}$ indices, we can rewrite the test in the simpler form

$$T_{\Omega} = \frac{\sum_{\mathbf{I}} \Omega_{\mathbf{I}}^{(n_{\mathbf{I}})} (\xi \dots \xi)_{\mathbf{I}}}{\sqrt{\text{Var} \left(\sum_{\mathbf{I}} \Omega_{\mathbf{I}}^{(n_{\mathbf{I}})} (\xi \dots \xi)_{\mathbf{I}} \right)}} \quad (58)$$

with the conditions

$$0 = \sum_{\mathbf{I}} \Omega_{\mathbf{I}}^{(n_{\mathbf{I}})} \mu_{\mathbf{I}}^{(m)} \quad , \quad m = 1, 2, 3 \dots \quad (59)$$

If we constrain these tests to be first order polynomials in ξ_i , we recover the linear case with $\Omega_i^{(1)} = \Omega_i / \xi_i^0$. Note that linear tests are always strongly centered. In fact in the infinite site model the spectrum is always proportional to θ , which consequently factorizes out by linearity and therefore has no effect on the centeredness.

⁶In other words, $\mu_{ijk\dots} = \sum_p \theta^p L^p \mu_{ijk\dots}^{(p)}$, where the coefficients $\mu_{ijk\dots}^{(p)}$ are independent on θ .

Weakly centered tests are tests that are centered but not strongly centered, i.e., they are centered if and only if $\hat{\theta} = \theta$. The general form for weakly centered tests is

$$T_{\Gamma} = \frac{\gamma + \sum_{i=1}^{n-1} \Gamma_i^{(1)} \xi_i + \sum_{i,j=1}^{n-1} \Gamma_{ij}^{(2)} \xi_i \xi_j + \sum_{i,j,k=1}^{n-1} \Gamma_{ijk}^{(3)} \xi_i \xi_j \xi_k + \dots}{\sqrt{\text{Var} \left(\gamma + \sum_{i=1}^{n-1} \Gamma_i^{(1)} \xi_i + \sum_{i,j=1}^{n-1} \Gamma_{ij}^{(2)} \xi_i \xi_j + \sum_{i,j,k=1}^{n-1} \Gamma_{ijk}^{(3)} \xi_i \xi_j \xi_k + \dots \right)}} \quad (60)$$

with the condition

$$0 = \gamma + \sum_{i=1}^{n-1} \Gamma_i^{(1)} \mu_i + \sum_{i,j=1}^{n-1} \Gamma_{ij}^{(2)} \mu_{ij} + \sum_{i,j,k=1}^{n-1} \Gamma_{ijk}^{(3)} \mu_{ijk} + \dots \quad (61)$$

where the $\Gamma_{ijk\dots}$ are real symmetric weights, possibly dependent on θ . We can simplify these expressions using the same notation as above, obtaining the simpler form

$$T_{\Gamma} = \frac{\gamma + \sum_{\mathbf{I}} \Gamma_{\mathbf{I}}^{(n_{\mathbf{I}})} (\xi \dots \xi)_{\mathbf{I}}}{\sqrt{\text{Var} \left(\gamma + \sum_{\mathbf{I}} \Gamma_{\mathbf{I}}^{(n_{\mathbf{I}})} (\xi \dots \xi)_{\mathbf{I}} \right)}} \quad (62)$$

with the condition

$$0 = \gamma + \sum_{\mathbf{I}} \Gamma_{\mathbf{I}}^{(n_{\mathbf{I}})} \mu_{\mathbf{I}} \quad (63)$$

Also for this class of tests, the sum can be limited to polynomials of fixed order or extended to power series. Note that the rescaled version of the G test by Fu presented above belongs to this class.

The important difference between strongly and weakly centered tests is related to the robustness with respect to a biased estimation of θ . Since the class of weakly centered tests is much larger than the class of strongly centered ones, it should be easier to find powerful tests in the former class than in the latter. However, even if weakly centered tests could be more powerful, they would not be centered in scenarios where the value of θ could not be estimated precisely. On the other side, strongly centered tests are robust with respect to a bad estimation of θ and therefore they would be preferable in scenarios where an unbiased estimation of θ is troublesome.

The scaling rule (15) can be generalized to nonlinear tests in terms of functions $\Omega^{(n_{\mathbf{I}})}(f_1, f_2 \dots f_{n_{\mathbf{I}}})$ for strongly centered and $\Gamma^{(n_{\mathbf{I}})}(f_1, f_2 \dots f_{n_{\mathbf{I}}})$ for

weakly centered tests:

$$\Omega_{\mathbf{I}}^{(n_{\mathbf{I}})} \simeq \frac{1}{n^{n_{\mathbf{I}}}} \Omega^{(n_{\mathbf{I}})} \left(\frac{i}{n}, \frac{j}{n}, \frac{k}{n} \dots \right) \quad (64)$$

$$\Gamma_{\mathbf{I}}^{(n_{\mathbf{I}})} \simeq \frac{1}{n^{n_{\mathbf{I}}}} \Gamma^{(n_{\mathbf{I}})} \left(\frac{i}{n}, \frac{j}{n}, \frac{k}{n} \dots \right) \quad (65)$$

Fixing the precise scaling is more ambiguous than in the linear case because there are many different ways to preserve centeredness. For this reason the choice of scaling would be different for strongly and weakly centered tests and will not be discussed here.

All the possible nonlinear neutrality tests based on the frequency spectrum fall into one of the two classes presented in this section and have the form (58),(59) or (62),(63). Since both these classes contain an infinite number of possible choices of weights, the only reasonable criterion to study general nonlinear tests is to select the most powerful or interesting ones. Apart from the Hotelling choice of Fu [11], the most interesting choice is apparently the subclass of nonlinear optimal tests, which will be discussed in the next sections.

5.2. Strongly centered optimal tests

As discussed for the linear case, optimal tests depend on the expected alternative scenario. In the nonlinear case, in principle it would be possible to find generic optimal tests, but there is no clear framework to obtain them. For this reason we limit our study to the case of optimal tests for a specific alternative scenario. We denote by $\bar{\mu}_{ijk\dots} = \mathcal{E}(\xi_i \xi_j \xi_k \dots)$ the moments of the alternative spectrum for this scenario.

Since we use the same normalization for linear and nonlinear tests, the optimality condition corresponds to the maximization of the expected value of the test under the alternative scenario

$$\mathcal{E}(T_{\Omega}) = \frac{\sum_{\mathbf{I}} \Omega_{\mathbf{I}}^{(n_{\mathbf{I}})} \bar{\mu}_{\mathbf{I}}}{\sqrt{\text{Var} \left(\sum_{\mathbf{I}} \Omega_{\mathbf{I}}^{(n_{\mathbf{I}})} (\xi \dots \xi)_{\mathbf{I}} \right)}} \quad (66)$$

and can be justified as in the linear case.

We denote by $\tilde{\mathbf{I}}$ the ordered sequence of the indices contained in $\mathbf{I} = ijk\dots$ and by $\sigma(\mathbf{I})$ the number of distinct permutations of the sequence \mathbf{I} , i.e. the

total number of permutations divided by the number of permutations that leave \mathbf{I} invariant. The main result for the optimal weights is presented in this theorem.

Theorem 6. *The maxima of $\mathcal{E}(T_\Omega)$ correspond to the weights*

$$\Omega_{\mathbf{I}}^{(n_{\mathbf{I}})} = \frac{1}{\sigma(\mathbf{I})} \left(\sum_{\tilde{\mathbf{I}}} C_{\tilde{\mathbf{I}}\tilde{\mathbf{L}}}^{-1} \bar{\mu}_{\tilde{\mathbf{L}}} - \sum_k \sum_l \sum_{\tilde{\mathbf{I}}} C_{\tilde{\mathbf{I}}\tilde{\mathbf{L}}}^{-1} \mu_{\tilde{\mathbf{L}}}^{(k)} \mathcal{M}_{kl} \sum_{\tilde{\mathbf{J}}, \tilde{\mathbf{K}}} \mu_{\tilde{\mathbf{J}}}^{(l)} C_{\tilde{\mathbf{J}}\tilde{\mathbf{K}}}^{-1} \bar{\mu}_{\tilde{\mathbf{K}}} \right) \quad (67)$$

where the matrices $C_{\tilde{\mathbf{I}}\tilde{\mathbf{J}}}^{-1}$ and \mathcal{M}_{kl} satisfy the identities

$$\sum_{\tilde{\mathbf{K}}} C_{\tilde{\mathbf{I}}\tilde{\mathbf{K}}}^{-1} (\mu_{\tilde{\mathbf{K}}\tilde{\mathbf{J}}} - \mu_{\tilde{\mathbf{K}}}\mu_{\tilde{\mathbf{J}}}) = \delta_{\tilde{\mathbf{I}}\tilde{\mathbf{J}}} \quad (68)$$

$$\sum_r \mathcal{M}_{kr} \sum_{\tilde{\mathbf{I}}, \tilde{\mathbf{L}}} \mu_{\tilde{\mathbf{I}}}^{(r)} C_{\tilde{\mathbf{I}}\tilde{\mathbf{L}}}^{-1} \mu_{\tilde{\mathbf{L}}}^{(l)} = \delta_{kl} \quad (69)$$

Moreover, the variance of the corresponding unnormalized test under the null model is equal to its expected value under the alternative model:

$$\text{Var} \left(\sum_{\mathbf{I}} \Omega_{\mathbf{I}}^{(n_{\mathbf{I}})} (\xi \dots \xi)_{\mathbf{I}} \right) = \sum_{\mathbf{I}} \Omega_{\mathbf{I}}^{(n_{\mathbf{I}})} \bar{\mu}_{\mathbf{I}} \quad (70)$$

Note that in general all the weights of the above optimal solution (67) are nonzero, therefore the maximum average value of the test for optimal tests built on polynomials of degree d increases with the degree d . This suggests that optimal tests of higher degree should be more powerful than linear optimal tests.

We provide explicit formulae for the above weights for the optimal quadratic test in the independent sites approximation. Given $E(\xi_i) = \mu_i$ and $\mathcal{E}(\xi_i) = \bar{\mu}_i$, the relevant weights $\Omega_{\mathbf{I}}^{(n_{\mathbf{I}})}$ are

$$\Omega_i^{(1)} = (\Sigma_\mu + 2 - \Sigma_{\bar{\mu}}) \left(\frac{\bar{\mu}_i}{\mu_i} - \frac{\Sigma_{\bar{\mu}}}{\Sigma_\mu} \right) - \frac{1}{2} \left(\frac{\bar{\mu}_i^2}{\mu_i^2} - \frac{\Sigma_{\bar{\mu}}^2}{\Sigma_\mu^2} \right) \quad (71)$$

$$\Omega_{ii}^{(2)} = - \left(\frac{\bar{\mu}_i}{\mu_i} - \frac{\Sigma_{\bar{\mu}}}{\Sigma_\mu} \right) + \frac{1}{2} \left(\frac{\bar{\mu}_i^2}{\mu_i^2} - \frac{\Sigma_{\bar{\mu}}^2}{\Sigma_\mu^2} \right) \quad (72)$$

$$\Omega_{ij}^{(2)} = \frac{1}{2} \left[\left(\frac{\bar{\mu}_i \bar{\mu}_j}{\mu_i \mu_j} - \frac{\Sigma_{\bar{\mu}}^2}{\Sigma_\mu^2} \right) - \left(\frac{\bar{\mu}_i}{\mu_i} - \frac{\Sigma_{\bar{\mu}}}{\Sigma_\mu} \right) - \left(\frac{\bar{\mu}_j}{\mu_j} - \frac{\Sigma_{\bar{\mu}}}{\Sigma_\mu} \right) \right] \quad (73)$$

where $\Sigma_\mu = \sum_{i=1}^{n-1} \mu_i$ and $\Sigma_{\bar{\mu}} = \sum_{i=1}^{n-1} \bar{\mu}_i$. All these formulae are also valid for the folded spectrum if the appropriate μ_i and $\bar{\mu}_i$ are used. These results are discussed in Appendix A.

For optimal tests of higher degree, explicit expressions become cumbersome and the numerical implementation of the test (67) and the matrices (68), (69) is more convenient.

5.3. Weakly centered optimal tests

In this case the optimality condition corresponds to the maximization of the expression

$$\mathcal{E}(T_\Gamma) = \frac{\gamma + \sum_{\mathbf{I}} \Gamma_{\mathbf{I}}^{(n_{\mathbf{I}})} \bar{\mu}_{\mathbf{I}}}{\sqrt{\text{Var} \left(\gamma + \sum_{\mathbf{I}} \Gamma_{\mathbf{I}}^{(n_{\mathbf{I}})} (\xi \dots \xi)_{\mathbf{I}} \right)}} \quad (74)$$

with the same condition

$$0 = \gamma + \sum_{\mathbf{I}} \Gamma_{\mathbf{I}}^{(n_{\mathbf{I}})} \mu_{\mathbf{I}} \quad (75)$$

The simplest case corresponds to a first order polynomial

$$T_\Gamma = \frac{\gamma + \sum_{i=1}^{n-1} \Gamma_i^{(1)} \xi_i}{\sqrt{\gamma^2 + 2\gamma \sum_{j=1}^{n-1} \Gamma_j^{(1)} \mu_j + \sum_{j=1}^{n-1} \sum_{k=1}^{n-1} \Gamma_j^{(1)} \Gamma_k^{(1)} \mu_{jk}}} \quad (76)$$

whose maximum corresponds to the optimal weights

$$\Gamma_i^{(1)} = \sum_{j=1}^{n-1} c_{ij}^{-1} (\bar{\mu}_j - \mu_j) \quad , \quad \gamma = - \sum_{j=1}^{n-1} \sum_{k=1}^{n-1} \mu_j c_{jk}^{-1} (\bar{\mu}_k - \mu_k) \quad (77)$$

where c_{ij}^{-1} is the inverse matrix of the covariance matrix $c_{ij} = \mu_{ij} - \mu_i \mu_j$. Since $\gamma \neq 0$ for this optimal test, the value of this test for the specific scenario for which it is built is larger than the value of the corresponding linear optimal test. In fact the maximum of the test is

$$\mathcal{E}(T_\Gamma) = \sqrt{\sum_{j=1}^{n-1} \sum_{k=1}^{n-1} (\bar{\mu}_j - \mu_j) c_{jk}^{-1} (\bar{\mu}_k - \mu_k)} \quad (78)$$

that should be compared to the maximum of the optimal test for the linear case, which can be rewritten as

$$\mathcal{E}(T_\Omega)_{linear} = \sqrt{\sum_{j=1}^{n-1} \sum_{k=1}^{n-1} (\bar{\mu}_j - \mu_j) c_{jk}^{-1} (\bar{\mu}_k - \mu_k) - \frac{\left(\sum_{j=1}^{n-1} \sum_{k=1}^{n-1} \mu_j c_{jk}^{-1} (\bar{\mu}_k - \mu_k)\right)^2}{\sum_{j=1}^{n-1} \sum_{k=1}^{n-1} \mu_j c_{jk}^{-1} \mu_k}} \quad (79)$$

The comparison shows clearly that nonlinear optimal tests are always more powerful than linear optimal tests for the same scenario.

The form of the results for the general case is similar to this simple case.

Theorem 7. *The maxima of $\mathcal{E}(T_\Gamma)$ correspond to the weights*

$$\Gamma_{\mathbf{I}}^{(n_{\mathbf{I}})} = \frac{1}{\sigma(\mathbf{I})} \sum_{\tilde{\mathbf{J}}} C_{\tilde{\mathbf{I}}\tilde{\mathbf{J}}}^{-1} (\bar{\mu}_{\tilde{\mathbf{J}}} - \mu_{\tilde{\mathbf{J}}}) \quad , \quad \gamma = - \sum_{\tilde{\mathbf{J}}} \mu_{\tilde{\mathbf{J}}} \sum_{\tilde{\mathbf{K}}} C_{\tilde{\mathbf{J}}\tilde{\mathbf{K}}}^{-1} (\bar{\mu}_{\tilde{\mathbf{K}}} - \mu_{\tilde{\mathbf{K}}}) \quad (80)$$

where $C_{\tilde{\mathbf{I}}\tilde{\mathbf{J}}}^{-1}$ satisfied the identity

$$\sum_{\tilde{\mathbf{K}}} C_{\tilde{\mathbf{I}}\tilde{\mathbf{K}}}^{-1} (\mu_{\tilde{\mathbf{K}}\tilde{\mathbf{J}}} - \mu_{\tilde{\mathbf{K}}} \mu_{\tilde{\mathbf{J}}}) = \delta_{\tilde{\mathbf{I}}\tilde{\mathbf{J}}} \quad (81)$$

Moreover, the variance of the corresponding unnormalized test under the null model is equal to its expected value under the alternative model:

$$\text{Var} \left(\gamma + \sum_{\mathbf{I}} \Gamma_{\mathbf{I}}^{(n_{\mathbf{I}})} (\xi \dots \xi)_{\mathbf{I}} \right) = \sum_{\mathbf{I}} \Gamma_{\mathbf{I}}^{(n_{\mathbf{I}})} \bar{\mu}_{\mathbf{I}} + \gamma \quad (82)$$

Also in this case, the power of optimal tests based on polynomials of higher degree increases with the degree of the polynomial.

It is possible to give explicit expressions of the above matrix and moments for the optimal quadratic test. The formulae for the weights $\Gamma_{\mathbf{I}}^{(n_{\mathbf{I}})}$ for the

unfolded spectrum are

$$\Gamma_i^{(1)} = (\Sigma_\mu + 2 - \Sigma_{\bar{\mu}}) \left(\frac{\bar{\mu}_i}{\mu_i} - 1 \right) - \frac{1}{2} \left(\frac{\bar{\mu}_i^2}{\mu_i^2} - 1 \right) \quad (83)$$

$$\Gamma_{ii}^{(2)} = \frac{1}{2} \left(\frac{\bar{\mu}_i}{\mu_i} - 1 \right)^2 \quad (84)$$

$$\Gamma_{ij}^{(2)} = \frac{1}{2} \left(\frac{\bar{\mu}_i}{\mu_i} - 1 \right) \left(\frac{\bar{\mu}_j}{\mu_j} - 1 \right) \quad (85)$$

$$\gamma = \frac{1}{2} (\Sigma_\mu - \Sigma_{\bar{\mu}}) (\Sigma_\mu + 2 - \Sigma_{\bar{\mu}}) \quad (86)$$

These results are valid in the independent sites approximation. They are also valid for the folded spectrum if the appropriate μ_i and $\bar{\mu}_i$ are used. An expression for the denominator of the test in the independent sites approximation can be found in Appendix A.

5.4. Simulations of the power of optimal tests

Since a theoretical evaluation of the power of optimal tests of different degree is not possible, we evaluate numerically the power of some of these tests in different scenarios. We consider the best possible case, that is, we assume that the precise value of θ is known. Moreover we assume unlinked sites and $\theta \ll 1$. In this approximation, as shown in Appendix A, the moments $E(\xi_i \xi_j \xi_k \dots)$ depend only on the first moments $\mu_i = \theta L \xi_i^0$ and similarly $\mathcal{E}(\xi_i \xi_j \xi_k \dots)$ depend only on $\bar{\mu}_i = \theta L \bar{\xi}_i$, therefore optimal tests depend only on the alternative and null average spectra.

Note that for numerical simulations of optimal tests of higher degree, the numerical implementation can be made easier if all the occurrences of inverse covariance matrices $C_{\mathbf{ij}}^{-1}$ in the the above formulae are replaced with the corresponding second moments $\mu_{\mathbf{ij}}^{-1}$, both in the expressions (67), (80) and in the definition (69). The test is the same because of the centeredness condition, as it can be verified explicitly.

We compare four optimal tests. The first two are the linear and quadratic strongly centered optimal tests, which are denoted by $T_{O(1)}^{sc}$ and $T_{O(2)}^{sc}$ respectively. The third test is the weakly centered optimal test $T_{O(1)}^{wc}$ based on a first order polynomial and presented in (76). The last optimal test $T_{O(2)}^{wc}$ is also weakly centered and based on on a quadratic polynomial. The explicit formulae for the computation of the weights of $T_{O(2)}^{sc}$ and $T_{O(2)}^{wc}$ were given in equations (71)-(73) and (83)-(86).

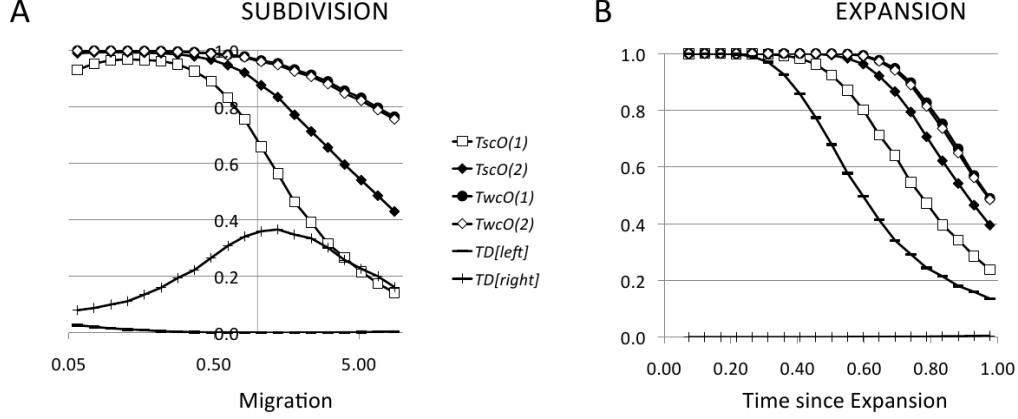


Figure 3: Statistical power of nonlinear optimal tests from coalescent simulations for the 5% tail, compared with Tajima’s D test (for the left and the right tail). The parameters for the simulations are: $n = 20$, $\theta = 0.05$, $L = 1000\text{bp}$, $\rho = \infty$; two populations considered but only one sampled (for panel A); expansion factor $N_0/N = 10$ (for panel B).

We simulated two demographic processes: (A) subdivision, where two populations having identical size exchange individuals given a symmetric migration rate M , then individuals are sampled from one population only; (B) expansion, where the population size changes by a factor $N_0/N = 10$ at a time T before present (in units of $4N$ generations). For each value of the parameters M and T , 10^6 simulations were performed with *mlcoalsim* v1.98b [20] for a region of 1000 bases with variability $\theta = 0.05$ and recombination $\rho = \infty$ and a sample size of $n = 20$ (haploid) individuals. Confidence intervals at 95% level were estimated from 10^6 simulations of the standard neutral coalescent with the same parameters.

In Figure 3 we compare the power of the tests in the best possible situation, namely when θ is known with good precision. In this condition all optimal tests should give the best results. In fact, the power of weakly centered tests ($T_{O(1)}^{wc}$ and $T_{O(2)}^{wc}$) is impressive, being around 100% for a large part of the parameter space and decreasing for large migration rates (Figure 3A) and long times (Figure 3B) as every other test, because the frequency spectrum for these cases becomes very similar to the standard spectrum. So,

weakly centered tests show a very good theoretical performance, counterbalanced by their lack of robustness. The power of $T_{O(1)}^{wc}$ and $T_{O(2)}^{wc}$ are almost identical, therefore the contribution of the quadratic part to $T_{O(2)}^{wc}$ is probably not relevant.

On the other hand, strongly centered optimal tests are quite more powerful than Tajima's D but less powerful than weakly centered tests, as expected. However, there is a sensible difference in power between $T_{O(1)}^{sc}$ and $T_{O(2)}^{sc}$: in the range of parameters where the power of weakly centered tests is around 100%, both strongly centered tests show a good performance not so far from the weakly centered ones, while in the less favourable range the quadratic test $T_{O(2)}^{sc}$, while performing worse than the weakly centered tests, has a power that is 20% higher than the linear test $T_{O(1)}^{sc}$. Taking into account the robustness of the tests, these simulations show that optimal tests like $T_{O(2)}^{sc}$ could be an interesting alternative to the usual linear tests.

6. Conclusions

In this paper we have presented a systematic analysis of neutrality tests based on the site frequency spectrum. This study is intended to extend and complete the recent works in [8] and [10] by extending the study of the linear neutrality tests recently presented by Achaz, their properties and the optimal tests that can be obtained in this framework; in a further generalization, we also consider the most general class of tests that can be written as a power series of the different components ξ_i of the spectrum. The aim of the paper is to give mathematical guidelines to build new and more effective tests. The proposed guidelines are the scaling relation (15) and the optimality condition based on the maximization of $\mathcal{E}(T_\Omega)$. Both these guidelines are thoroughly explained and discussed.

While nonlinear optimal tests have been shown to be more powerful than linear ones (and weakly centered tests more powerful than strongly centered ones), power is not the only important issue: also robustness must be taken into consideration. In fact there are three important remarks on the relative robustness of these tests. The first one is that, as already discussed, centeredness of weakly centered tests is not robust with respect to a biased estimate of θ , therefore these tests should be preferred to strongly centered tests only in situations where the value of θ is well known or a good estimate is available.

The second remark is that neither the weights nor the results of linear optimal tests do depend on the value of θ and on the number of segregating sites S , while the weights of nonlinear optimal tests have an explicit dependence on θ and their results depend not only on the spectrum but also on S , therefore the interpretation of the results of these tests is more complicated. However, this is not necessarily true for homogeneous tests of any degree, like the quadratic G_ξ test by Fu. An interesting development of this work could be a study of homogeneous tests of a given degree k satisfying the optimality condition, which can be easily obtained from equations (67) by restricting all ordered sequences of indices $\tilde{\mathbf{I}}, \tilde{\mathbf{J}}, \tilde{\mathbf{K}}, \tilde{\mathbf{L}}$ to contain precisely k indices (along with some “traceless” condition, in case). These homogeneous optimal tests (or at least some subclass of them) should depend weakly on S .

The third remark is that linear optimal tests have two interesting properties that are not shared by other tests: they depend only on the deviations from the null spectrum and they have an easy interpretation in terms of these deviations, that is, they are positive if the observed deviations are similar to the expected ones and negative if the observed deviations are opposite to the expected ones. These features give an important advantage to linear optimal tests.

Tests based on the frequency spectrum of polymorphic sites are fast, being based on simple matrix multiplications, and can be therefore applied to genome-wide data. Moreover, they can be used as summary statistics for Approximate Bayesian Computation or other statistical approaches to the analysis of sequence data. While linear tests are often used in this way, the nonlinear tests presented in this paper contain more information (related to the covariances and higher moments of the frequency spectrum) that could increase the power of these analyses.

Acknowledgments:

We thank M. Perez-Enciso and J. Rozas for useful comments. G.M. acknowledges support by Fondazione A. Della Riccia and Japan Society for Promotion of Science. Work funded by grant CGL2009-0934 (MICINN, Spain) to S.R.O.

A. Moments of the frequency spectrum in the independent sites approximation

We consider the limit $\theta \ll 1$, $L \rightarrow \infty$ with constant θL . The spectrum ξ_i can be written as a sum of spectra for all sites

$$\xi_i = \sum_{s=1}^L \xi_i(s) \quad (87)$$

where each variable $\xi_i(s)$ has a Bernoulli distribution $\xi_i(s) \in \{0, 1\}$ with probabilities $p(1) = \theta \xi_i^0$ and $p(0) = 1 - \theta \xi_i^0$ where $\xi_i^0 = 1/i$ under the standard neutral model. The expectation value of ξ_i is therefore

$$E(\xi_i) = \mu_i = \sum_{s=1}^L E(\xi_i(s)) = \theta L \xi_i^0 = \frac{\theta L}{i} \quad (88)$$

and similarly $\mathcal{E}(\xi_i) = \bar{\mu}_i = \theta L \bar{\xi}_i$ for a general model with average spectrum $\bar{\xi}_i$.

In the independent sites approximation, which is equivalent to the infinite recombination limit, the variables $\xi_i(s)$ and $\xi_i(s')$ are i.i.d. random variables for $s \neq s'$, and more generally the random variables $\xi_i(s)$ and $\xi_j(s')$ are independent for $s \neq s'$. The moments for a single site s can be calculated as

$$\begin{aligned} E(\xi_i(s)\xi_j(s)\xi_k(s)\dots) &= \sum_{a,b,c,\dots \in \{0,1\}} abc\dots P(\xi_i(s)=a, \xi_j(s)=b, \xi_k(s)=c, \dots) = \\ &= P(\xi_i(s)=1, \xi_j(s)=1, \xi_k(s)=1, \dots) = P(\xi_i(s)=1)\delta_{ij}\delta_{jk}\dots \end{aligned} \quad (89)$$

because the sum $\sum_{i=1}^{n-1} \xi_i(s) \in \{0, 1\}$, that is, different allele frequencies for the same site are mutually exclusive. Therefore all these moments are always linear in θ and are nonzero only when all indices are equal:

$$E(\xi_i(s)\xi_j(s)\xi_k(s)\dots) = \theta \xi_i^0 \delta_{ij} \delta_{jk} \dots = \frac{\theta}{i} \delta_{ij} \delta_{jk} \dots \quad (90)$$

Then the second moment can be evaluated as

$$\begin{aligned} E(\xi_i \xi_j) &= \sum_{s,s'} E(\xi_i(s)\xi_j(s')) = \sum_s E(\xi_i(s)\xi_j(s)) + \sum_{s \neq s'} E(\xi_i(s))E(\xi_j(s')) = \\ &= L \delta_{ij} \frac{\theta}{i} + L(L-1) \frac{\theta^2}{i j} = \delta_{ij} \mu_i + \mu_i \mu_j - \frac{\mu_i \mu_j}{L} \end{aligned} \quad (91)$$

and neglecting subleading orders in θ or L^{-1} like the last term above, we can calculate the third and forth moments that are needed for the calculation of C^{-1} .

The final result for the moments of the spectrum ξ_i is

$$E(\xi_i \xi_j) = \mu_{ij} = \delta_{ij} \mu_i + \mu_i \mu_j \quad (92)$$

$$E(\xi_i \xi_j \xi_k) = \mu_{ijk} = \delta_{ij} \delta_{jk} \mu_i + (\delta_{ik} \mu_i \mu_j + \delta_{jk} \mu_i \mu_j + \delta_{ij} \mu_i \mu_k) + \mu_i \mu_j \mu_k \quad (93)$$

$$\begin{aligned} E(\xi_i \xi_j \xi_k \xi_l) = \mu_{ijkl} = & \delta_{ij} \delta_{jk} \delta_{kl} \mu_i + (\delta_{ik} \delta_{jl} \mu_i \mu_j + \delta_{il} \delta_{jk} \mu_i \mu_j + \delta_{ij} \delta_{kl} \mu_i \mu_k) + \\ & + (\delta_{ij} \delta_{jk} \mu_i \mu_l + \delta_{ij} \delta_{jl} \mu_i \mu_k + \delta_{ik} \delta_{kl} \mu_i \mu_j + \delta_{jk} \delta_{kl} \mu_i \mu_j) + \\ & + (\delta_{il} \mu_i \mu_j \mu_k + \delta_{jl} \mu_i \mu_j \mu_k + \delta_{ik} \mu_i \mu_j \mu_l + \delta_{jk} \mu_i \mu_j \mu_l + \\ & + \delta_{ij} \mu_i \mu_k \mu_l + \delta_{kl} \mu_i \mu_j \mu_k) + \mu_i \mu_j \mu_k \mu_l \end{aligned} \quad (94)$$

All the results above can be applied to a general model simply by substituting μ_i with $\bar{\mu}_i$. Moreover they can be applied to the folded spectrum by taking $\mu_i = \theta L n / i(n-i)(1 + \delta_{n,2i})$ for the standard neutral model or $\bar{\mu}_i = \theta L (\xi_i + \bar{\xi}_{n-i}) / (1 + \delta_{n,2i})$ for general models.

We define some quantities in order to simplify the expressions for the weights:

$$\Sigma_\mu = \sum_{i=1}^{n-1} \mu_i \quad , \quad \Sigma_{\bar{\mu}} = \sum_{i=1}^{n-1} \bar{\mu}_i \quad , \quad \Sigma_q = \sum_{i=1}^{n-1} \frac{\bar{\mu}_i^2}{\mu_i} \quad (95)$$

If the spectrum is folded, all the sums in the above expressions run from 1 to $\lfloor n/2 \rfloor$.

The covariance matrix $C_{\mathbf{I}, \mathbf{J}}$ is

$$C_{i,j} = \mu_{ij} - \mu_i \mu_j = \delta_{ij} \mu_i \quad (96)$$

$$C_{ij,k} = \mu_{ijk} - \mu_{ij} \mu_k = \delta_{ij} \delta_{jk} \mu_i + (\delta_{ik} \mu_i \mu_j + \delta_{jk} \mu_i \mu_j) \quad (97)$$

$$\begin{aligned} C_{ij,kl} = \mu_{ijkl} - \mu_{ij} \mu_{kl} = & \delta_{ij} \delta_{jk} \delta_{kl} \mu_i + (\delta_{ik} \delta_{jl} \mu_i \mu_j + \delta_{il} \delta_{jk} \mu_i \mu_j) + \\ & + (\delta_{ij} \delta_{jk} \mu_i \mu_l + \delta_{ij} \delta_{jl} \mu_i \mu_k + \delta_{ik} \delta_{kl} \mu_i \mu_j + \delta_{jk} \delta_{kl} \mu_i \mu_j) + \\ & + (\delta_{il} \mu_i \mu_j \mu_k + \delta_{jl} \mu_i \mu_j \mu_k + \delta_{ik} \mu_i \mu_j \mu_l + \delta_{jk} \mu_i \mu_j \mu_l) \end{aligned} \quad (98)$$

with the elements $C_{ij,k}$ and $C_{ij,kl}$ that should be considered only for $i \leq j$,

$k \leq l$. It can be verified that the inverse matrix $C_{\tilde{\mathbf{I}}, \tilde{\mathbf{J}}}^{-1}$ has the form

$$C_{i,j}^{-1} = 1 + \delta_{ij} \frac{2\mu_i (\Sigma_\mu + 3) + 1}{2\mu_i^2} \quad (99)$$

$$C_{ij,k}^{-1} = \delta_{ij} \delta_{jk} \frac{2\mu_i - 1}{2\mu_i^2} - (\delta_{ik} + \delta_{jk}) \frac{1}{\mu_k} \quad (100)$$

$$C_{ij,kl}^{-1} = (\delta_{ik} \delta_{jl} + \delta_{il} \delta_{jk}) \frac{1}{\mu_i \mu_j} - \frac{3}{2} \delta_{ij} \delta_{ik} \delta_{jl} \frac{1}{\mu_i^2} \quad (101)$$

The formulae for the weakly centered quadratic test (83-86) can be obtained from these formulae and the definition (80). The corresponding variance in the denominator (62) is

$$\begin{aligned} \text{Var} \left(\sum_{\mathbf{I}} \Gamma_{\mathbf{I}}^{(n_{\mathbf{I}})} (\xi \dots \xi)_{\mathbf{I}} \right) &= 2(\Sigma_{\bar{\mu}} - \Sigma_q/2)^2 + \Sigma_\mu (\Sigma_\mu/2 + 1 - 2\Sigma_{\bar{\mu}} + \Sigma_q) \\ &\quad + \Sigma_q - 2\Sigma_{\bar{\mu}} \end{aligned} \quad (102)$$

The matrix \mathcal{M} is the inverse of the matrix $\mathcal{M}_{rl}^{-1} = \sum_{\tilde{\mathbf{I}}, \tilde{\mathbf{L}}} \mu_{\tilde{\mathbf{I}}}^{(r)} C_{\tilde{\mathbf{I}}\tilde{\mathbf{L}}}^{-1} \mu_{\tilde{\mathbf{L}}}^{(l)}$, which can be easily calculated from the above equations as

$$\mathcal{M}_{11}^{-1} = 2\Sigma_\mu^2 + \Sigma_\mu \quad (103)$$

$$\mathcal{M}_{12}^{-1} = -\Sigma_\mu^2 \quad (104)$$

$$\mathcal{M}_{22}^{-1} = \Sigma_\mu^2/2 \quad (105)$$

and the matrix \mathcal{M} is

$$\mathcal{M} = \begin{pmatrix} 1/\Sigma_\mu & 2/\Sigma_\mu \\ 2/\Sigma_\mu & 4/\Sigma_\mu + 2/\Sigma_\mu^2 \end{pmatrix} \quad (106)$$

The formulae (71-73) can be obtained from the formula (106) and from the following results:

$$\sum_{\tilde{\mathbf{I}}} C_{i,\tilde{\mathbf{I}}}^{-1} \bar{\mu}_{\tilde{\mathbf{I}}} = \Sigma_{\bar{\mu}} + (\Sigma_\mu + 2 - \Sigma_{\bar{\mu}}) \frac{\bar{\mu}_i}{\mu_i} - \frac{1}{2} \left(\frac{\bar{\mu}_i}{\mu_i} \right)^2 \quad (107)$$

$$\sum_{\tilde{\mathbf{I}}} C_{ii,\tilde{\mathbf{I}}}^{-1} \bar{\mu}_{\tilde{\mathbf{I}}} = -\frac{\bar{\mu}_i}{\mu_i} + \frac{1}{2} \left(\frac{\bar{\mu}_i}{\mu_i} \right)^2 \quad (108)$$

$$\sum_{\tilde{\mathbf{I}}} C_{ij,\tilde{\mathbf{I}}}^{-1} \bar{\mu}_{\tilde{\mathbf{I}}} = -\frac{\bar{\mu}_i}{\mu_i} - \frac{\bar{\mu}_j}{\mu_j} + \frac{\bar{\mu}_i \bar{\mu}_j}{\mu_i \mu_j} \quad (109)$$

$$\sum_{\tilde{\mathbf{I}}, \tilde{\mathbf{J}}} \mu_{\tilde{\mathbf{I}}}^{(1)} C_{\tilde{\mathbf{I}}\tilde{\mathbf{J}}}^{-1} \bar{\mu}_{\tilde{\mathbf{J}}} = \Sigma_{\bar{\mu}}(2\Sigma_{\mu} - \Sigma_{\bar{\mu}} + 1) \quad \sum_{\tilde{\mathbf{I}}, \tilde{\mathbf{J}}} \mu_{\tilde{\mathbf{I}}}^{(2)} C_{\tilde{\mathbf{I}}\tilde{\mathbf{J}}}^{-1} \bar{\mu}_{\tilde{\mathbf{J}}} = \Sigma_{\bar{\mu}}^2/2 - \Sigma_{\mu}\Sigma_{\bar{\mu}} \quad (110)$$

$$\sum_{\tilde{\mathbf{I}}} C_{i, \tilde{\mathbf{I}}}^{-1} \mu_{\tilde{\mathbf{I}}}^{(1)} = 2\Sigma_{\mu} + 2 \quad \sum_{\tilde{\mathbf{I}}} C_{i, \tilde{\mathbf{I}}}^{-1} \mu_{\tilde{\mathbf{I}}}^{(2)} = -\Sigma_{\mu} - 1/2 \quad (111)$$

$$\sum_{\tilde{\mathbf{I}}} C_{ii, \tilde{\mathbf{I}}}^{-1} \mu_{\tilde{\mathbf{I}}}^{(1)} = -1 \quad \sum_{\tilde{\mathbf{I}}} C_{ii, \tilde{\mathbf{I}}}^{-1} \mu_{\tilde{\mathbf{I}}}^{(2)} = 1/2 \quad (112)$$

$$\sum_{\tilde{\mathbf{I}}} C_{ij, \tilde{\mathbf{I}}}^{-1} \mu_{\tilde{\mathbf{I}}}^{(1)} = -2 \quad \sum_{\tilde{\mathbf{I}}} C_{ij, \tilde{\mathbf{I}}}^{-1} \mu_{\tilde{\mathbf{I}}}^{(2)} = 1 \quad (113)$$

B. Proofs

PROOF OF THEOREM 3. Choose a vector Δ_i in \mathbb{R}^{n-1} that is orthogonal both to $i\Omega_i$ and to i , that is, such that $\sum_i i\Delta_i\Omega_i = 0$ and $\sum_i i\Omega_i = 0$. Since $\alpha/ia_n + (1 - \alpha)\Delta_i$ is a set of continuous functions of α , its minimum is also continuous in α . Moreover, the minimum is clearly positive if $\alpha = 1$, while it is negative by construction if $\alpha = 0$. The theorem follows from the intermediate value theorem. ■

PROOF OF THEOREM 4. The proof is similar to the previous one. Choose a function $\Delta(f)$ in $C^\infty(0, 1)$ to satisfy both $\int_{1/N}^1 df f\Omega(f)\Delta(f) = 0$ and $\int_{1/N}^1 df f\Delta(f) = 0$. (Since these conditions correspond just to two independent functionals of $\Delta(f)$ and $C^\infty(0, 1)$ is an infinite-dimensional linear space, the existence of such a function is guaranteed.) Since $\alpha/f \log(N) + (1 - \alpha)\Delta(f)$ is a continuous functions of α and its infimum is not $\pm\infty$, its infimum is also continuous in α . Moreover, the infimum is clearly positive if $\alpha = 1$, while it is negative by construction if $\alpha = 0$. The theorem follows from the intermediate value theorem. ■

PROOF OF THEOREM 5. The vectors Ω_i are a basis of the subspace $\mathbb{R}^{n-2} \subset \mathbb{R}^{n-1}$ defined by the condition $\sum_i \Omega_i = 0$, that is, the space of vectors orthogonal to the vector whose components are $v_i = 1$. Therefore the only vectors $i\Delta_i$ that are orthogonal to all the vectors in this basis are precisely of the form $i\Delta_i \propto v_i$, that is, $\Delta_i = \text{const}/i$. ■

The theorems on the form of optimal tests can be easily proved from a general lemma.

Lemma 1. Consider a function $f : \mathbb{R}^M \setminus \{0\} \rightarrow \mathbb{R}$ of the form

$$f(\vec{v}) = \frac{\vec{v} \cdot \vec{w}}{\sqrt{\vec{v} \cdot Q \vec{v}}} \quad (114)$$

where $\vec{w} \in \mathbb{R}^M$ and Q is a $M \times M$ symmetric positive matrix, and a $K \times M$ matrix R with $K < M$ and maximum rank. The extrema of the function f restricted to the subspace $R\vec{v} = 0$ are given by

$$\vec{v}_\alpha = \alpha \left(Q^{-1} \vec{w} - Q^{-1} R^t (R Q^{-1} R^t)^{-1} R Q^{-1} \vec{w} \right) \quad (115)$$

The extrema with $\alpha > 0$ are maxima and the extrema with $\alpha < 0$ are minima of the function f . These extrema satisfy the identity

$$\vec{v}_\alpha \cdot Q \vec{v}_\alpha = \alpha \vec{v}_\alpha \cdot \vec{w} \quad (116)$$

PROOF. The existence of maxima and minima can be proved by the Weierstrass extreme value theorem. In fact f is continuous and invariant under a homothety with center in the origin of \mathbb{R}^M and positive scale factor, therefore the codomain of the function on the linear subspace defined by $R\vec{v} = 0$ is the same as the codomain of its restriction to the submanifold of unit vectors $|\vec{v}| = 1$, which is a compact space. The restriction of f is also continuous and the conclusion follows. To determine the extrema, the method of Lagrange multipliers states that it is sufficient to extremize the function

$$F(\vec{v}, \vec{\lambda}) = \frac{\vec{v} \cdot \vec{w}}{\sqrt{\vec{v} \cdot Q \vec{v}}} + \vec{\lambda} \cdot R \vec{v} \quad (117)$$

and since there are no boundaries, this is equivalent to the solution of the equations

$$0 = \vec{\nabla}_v f = \frac{\vec{w}}{\sqrt{\vec{v} \cdot Q \vec{v}}} - \frac{\vec{v} \cdot \vec{w}}{(\vec{v} \cdot Q \vec{v})^{3/2}} Q \vec{v} + R^t \vec{\lambda} \quad (118)$$

$$0 = \vec{\nabla}_\lambda f = R \vec{v} \quad (119)$$

The solution satisfies

$$\frac{\vec{v} \cdot \vec{w}}{\vec{v} \cdot Q \vec{v}} \vec{v} = Q^{-1} \vec{w} + Q^{-1} R^t \vec{\lambda} \quad (120)$$

and multiplying it by R and using (119) we obtain

$$\vec{\lambda} = - (RQ^{-1}R^t)^{-1} RQ^{-1}\vec{w} + \vec{l} \quad , \quad R^t\vec{l} = 0 \quad (121)$$

that can be inserted again in equation (120) to eliminate $\vec{\lambda}$. The resulting equation in \vec{v} admits only solutions of the form (115) and by substituting (115) into it, it can be checked that all values of $\alpha \neq 0$ correspond to solutions of (118),(119). The invariance of f under a central homothety with positive scale factor implies that the value of the function does not depend on $|\alpha|$. The function is positive for $\alpha > 0$ and negative for $\alpha < 0$, therefore solutions with $\alpha > 0$ correspond to maxima and solutions with $\alpha < 0$ correspond to minima. The identity (116) can be proved by substituting the solution (115). ■

PROOF OF THEOREM 6. The expected values of the tests (58) have the same functional form as the function f of Lemma 1. The correspondence is the following:

$$\vec{v} \rightarrow v_{\tilde{\mathbf{I}}} = \sigma(\tilde{\mathbf{I}})\Omega_{\tilde{\mathbf{I}}}^{(n_{\tilde{\mathbf{I}}})} \quad (122)$$

$$\vec{w} \rightarrow w_{\tilde{\mathbf{I}}} = \bar{\mu}_{\tilde{\mathbf{I}}} \quad (123)$$

$$Q \rightarrow Q_{\tilde{\mathbf{I}}\tilde{\mathbf{J}}} = \mu_{\tilde{\mathbf{I}}\tilde{\mathbf{J}}} - \mu_{\tilde{\mathbf{I}}}\mu_{\tilde{\mathbf{J}}} \quad (124)$$

$$R \rightarrow R_{k,\tilde{\mathbf{I}}} = \mu_{\tilde{\mathbf{I}}}^{(k)} \quad (125)$$

and the positivity of the matrix Q is guaranteed by the positivity of the variance for all possible choices of the weights. Application of Lemma 1 with $\alpha\alpha = 1$ gives the result (67). ■

PROOF OF THEOREM 7. We can immediately solve equation (75) for γ and substitute it in equation (74). Then γ is a function of the other weights and the maximization is unconstrained. It can be seen that also in this case, the expected values of the tests have the same functional form as the function f of Lemma 1. The correspondence is the following:

$$\vec{v} \rightarrow v_{\tilde{\mathbf{I}}} = \sigma(\tilde{\mathbf{I}})\Gamma_{\tilde{\mathbf{I}}}^{(n_{\tilde{\mathbf{I}}})} \quad (126)$$

$$\vec{w} \rightarrow w_{\tilde{\mathbf{I}}} = \bar{\mu}_{\tilde{\mathbf{I}}} - \mu_{\tilde{\mathbf{I}}} \quad (127)$$

$$Q \rightarrow Q_{\tilde{\mathbf{I}}\tilde{\mathbf{J}}} = \mu_{\tilde{\mathbf{I}}\tilde{\mathbf{J}}} - \mu_{\tilde{\mathbf{I}}}\mu_{\tilde{\mathbf{J}}} \quad (128)$$

$$R \rightarrow \text{empty } 0 \times M \text{ matrix} \quad (129)$$

and the positivity of the matrix Q is implied by the positivity of the variance. Then the result (80) follows from Lemma 1 with $\alpha = 1$. ■

References

- [1] M. Kreitman, Nucleotide polymorphism at the alcohol dehydrogenase locus of *Drosophila melanogaster*, *Nature* 304 (5925) (1983) 412–417.
- [2] R. Hudson, M. Kreitman, M. Aguadé, A test of neutral molecular evolution based on nucleotide data, *Genetics* 116 (1) (1987) 153.
- [3] R. Lewontin, J. Krakauer, Distribution of gene frequency as a test of the theory of the selective neutrality of polymorphisms, *Genetics* 74 (1) (1973) 175.
- [4] P. Sabeti, D. Reich, J. Higgins, H. Levine, D. Richter, S. Schaffner, S. Gabriel, J. Platko, N. Patterson, G. McDonald, et al., Detecting recent positive selection in the human genome from haplotype structure, *Nature* 419 (6909) (2002) 832–837.
- [5] F. Tajima, Statistical method for testing the neutral mutation hypothesis by DNA polymorphism, *Genetics* 123 (3) (1989) 585.
- [6] Y. Fu, W. Li, Statistical tests of neutrality of mutations, *Genetics* 133 (3) (1993) 693.
- [7] J. Fay, C. Wu, Hitchhiking under positive Darwinian selection, *Genetics* 155 (3) (2000) 1405.
- [8] G. Achaz, Frequency Spectrum Neutrality Tests: One for All and All for One, *Genetics* 183 (1) (2009) 249.
- [9] Y. Fu, Statistical properties of segregating sites, *Theoretical Population Biology* 48 (2) (1995) 172–197.
- [10] L. Ferretti, M. Perez-Enciso, S. Ramos-Onsins, Optimal Neutrality Tests Based on Frequency Spectrum, *Genetics* (2010) genetics.110.118570.
- [11] Y. Fu, New statistical tests of neutrality for DNA samples from a population, *Genetics* 143 (1) (1996) 557.
- [12] S. Schaeffer, Molecular population genetics of sequence length diversity in the *Adh* region of *Drosophila pseudoobscura*, *Genetics Research* 80 (03) (2003) 163–175.

- [13] K. Schmid, S. Ramos-Onsins, H. Ringys-Beckstein, B. Weisshaar, T. Mitchell-Olds, A multilocus sequence survey in *Arabidopsis thaliana* reveals a genome-wide departure from a neutral model of DNA sequence polymorphism, *Genetics* 169 (3) (2005) 1601.
- [14] S. Hutter, H. Li, S. Beisswanger, D. De Lorenzo, W. Stephan, Distinctly different sex ratios in African and European populations of *Drosophila melanogaster* inferred from chromosomewide single nucleotide polymorphism data, *Genetics* 177 (1) (2007) 469.
- [15] Y. Fu, Statistical tests of neutrality of mutations against population growth, hitchhiking and background selection, *Genetics* 147 (2) (1997) 915.
- [16] K. Zeng, Y. Fu, S. Shi, C. Wu, Statistical tests for detecting positive selection by utilizing high-frequency variants, *Genetics* 174 (3) (2006) 1431.
- [17] F. Tajima, Evolutionary relationship of DNA sequences in finite populations, *Genetics* 105 (2) (1983) 437.
- [18] G. Watterson, On the number of segregating sites in genetical models without recombination., *Theoretical population biology* 7 (2) (1975) 256.
- [19] G. Achaz, Testing for neutrality in samples with sequencing errors, *Genetics* 179 (3) (2008) 1409.
- [20] S. E. Ramos-Onsins, T. Mitchell-Olds, Mlcoalsim: multilocus coalescent simulations., *Evol Bioinform Online* 3 (2007) 41–44.